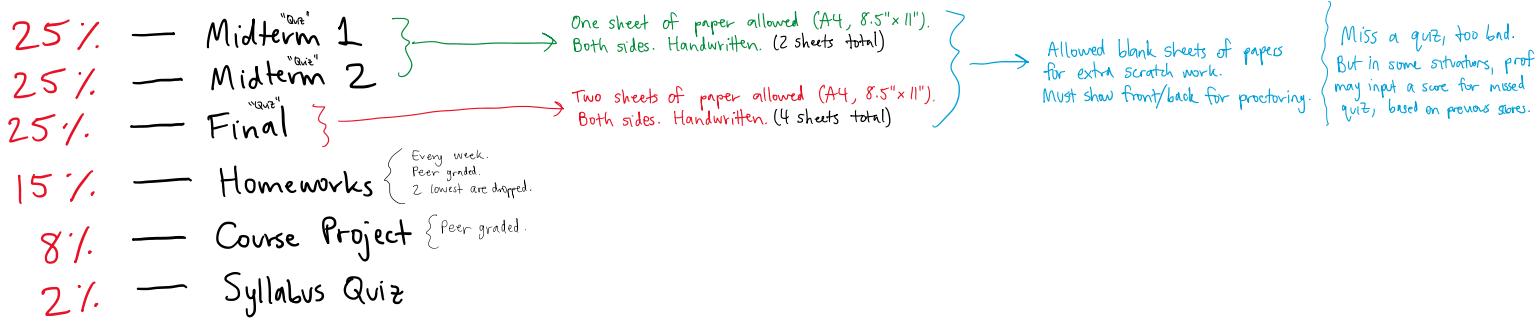


Syllabus

Sunday, January 9, 2022 3:37 PM



- Must finish:
- Onboarding Quiz (ungraded) — Due Jan 24, 2 AM
 - Syllabus Quiz (2%, unlimited tries) — Due Jan 24, 2 AM
 - Piazza Registration (part of onboarding)

Table of Contents

Friday, March 4, 2022 3:44 PM

- Week 1: SVM, kNN, Scaling
2: Validation, Clustering
3: Outliers, CUSUM
4: Exponential Smoothing, ARIMA, GARCH
5: Simple Linear Regression, Likelihood
6: Box-Cox Transformation, Detrending, PCA
7: CART, Confusion Matrices, Advanced Regression
8: Variable Selection
9: Design of Experiments, Probability, Queuing
10: Missing Data, Optimization
11: (In-class optimization exercises)
12: Advanced Optimization, Advanced Miscellaneous Modeling

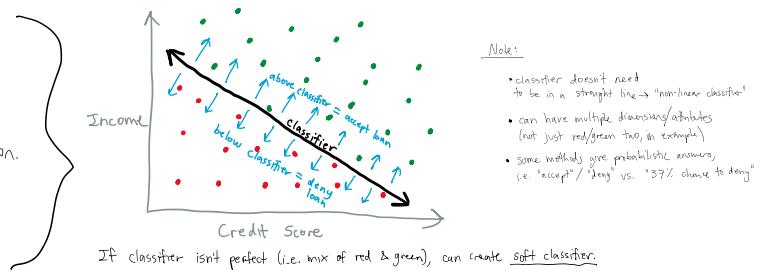
Note

characteristic/measurement of data → column of data
= attribute
= feature
= predictor
= covariate

Classification: the separation of data into two or more categories, or (a point's classification) the category a data point is put into

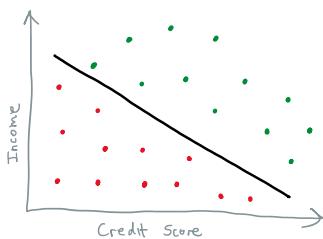
Classifier: A boundary that separates the data in two or more categories. Also (more generally) an algorithm that performs classification.

Support Vector Machines (SVM): classification algorithm that uses a boundary to separate the data into two or more categories ("classes")



↳ basically algorithm (math) to draw lines between points.

SVM Mathematics



$$n = \text{number of data points}$$

$$m = \text{number of attributes}$$

	Eye Color	Hair Color
1	Brown	Black
2	Black	Black
3	Blue	Blonde

$n=3$ (rows)

$m=2$ (Eye color, hair color)

$$x_{ij} = \text{ith attribute of jth data point}$$

$$x_{1j} = \text{credit score of person } j$$

$$x_{2j} = \text{income of person } j$$

$$y_j = \text{response for data point } j$$

$$\text{i.e. } y_j = \begin{cases} 1, & \text{if data point } j \text{ is green} \\ -1, & \text{if data point } j \text{ is red} \end{cases}$$

arbitrarily set as ± 1 ; can be any \pm constant

	Credit Score	Income
person 1	$600 = x_{11}$	$50,000 = x_{21}$
person 2	$650 = x_{12}$	$55,000 = x_{22}$
person 3	$700 = x_{13}$	$60,000 = x_{23}$
person 4		

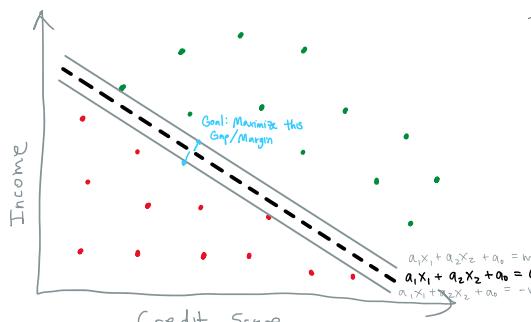
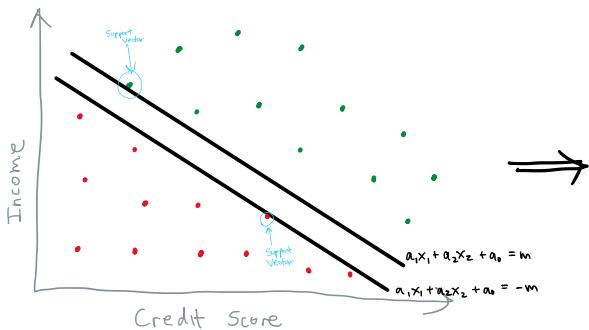
$$\text{Line: } a_1x_1 + a_2x_2 + \dots + a_mx_m + a_0 = 0 \quad (\text{from generic distance of a line equation})$$

$$\text{In other words: } \sum_{i=1}^m a_i x_{ij} + a_0 = 0$$

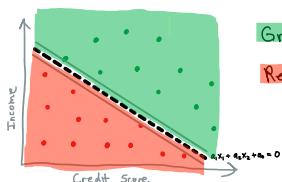
$a_i = \text{coefficient defining classifier}$ like a "summary" of n data points (row)
 $a_0 = \text{intercept of classifier}$ for every m attribute (column) vs. one-by-one

Note: near-zero a_i coefficients are likely not relevant for classification.

Can draw 2 parallel lines so a_0 is the classifier in the middle.



More Equations



$$\text{Green Points: } a_1x_{1j} + a_2x_{2j} + \dots + a_mx_{mj} + a_0 \geq 1$$

$$\text{Red Points: } a_1x_{1j} + a_2x_{2j} + \dots + a_mx_{mj} + a_0 \leq -1$$

could be any m number... picked 1 for scaling

$$\text{ALL POINTS: } (a_1x_{1j} + a_2x_{2j} + \dots + a_mx_{mj} + a_0)y_j \geq 1$$

$$\text{Distance between two solid lines} = \frac{2}{\sqrt{\sum_i (a_i)^2}}$$

(from linear algebra) Remember the goal is to maximize the margin. $\Rightarrow \text{distance} = \frac{(\text{constant})}{\sqrt{\text{deno}}} \Rightarrow \text{minimize } a_i! \Rightarrow \min_i \sum (a_i)^2$

$$\Rightarrow \text{Minimize } \sum_{i=1}^m (a_i)^2, \text{ subject to } (a_1x_{1j} + a_2x_{2j} + \dots + a_mx_{mj} + a_0)y_j \geq 1 \text{ for each data point.}$$

So, need to find values a_0 thru a_m that will maximize the margin, but can only choose values of coefficient a that correctly separate all the points.... software does this!

Summary:

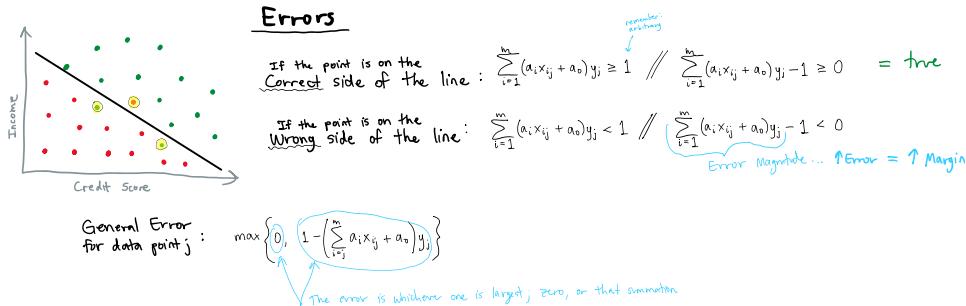
Given data $\{x_{ij}\}$ where x_{ij} is the j th factor of data point i :

- Classifier is the line: $0 = a_0 + \sum_j a_j x_j$
- Maximize gap between lines by minimizing $\sum_j a_j^2$

Soft Classifier

Use if not possible to perfectly separate red and green points.

Need to account for errors in classification, and trade off minimizing the errors we make vs. maximizing the margin.



Trade-Off: ↑ margins vs. ↓ Errors

Sometimes, having some error will increase our margin.
Other times, too many mistakes will not result in an inconclusive margin.

How do we calculate this trade-off?

This is called the Support Vector Machine (SVM) Model:

$$\underset{a_0, \dots, a_m}{\text{minimize}} \sum_{j=1}^n \max\left\{0, 1 - \left(\sum_{i=1}^m a_i x_{ij} + a_0\right) y_j\right\} + \lambda \sum_{i=1}^m (a_i)^2$$

$$\underset{a_0, \dots, a_m}{\text{minimize}} \sum_{j=1}^n \max\left\{0, 1 - \left(\sum_{i=1}^m a_i x_{ij} + a_0\right) y_j\right\} + \lambda \sum_{i=1}^m (a_i)^2$$

↑ Error
↓ Margin

Pick a value:
 $\uparrow \lambda = \uparrow \text{Margin (towards } \infty\text{)} \rightarrow \uparrow \text{error} \rightarrow \text{a large margin outweighs avoiding mistakes and classifying known data points.}$
 $\downarrow \lambda = \downarrow \text{Margin (towards } 0\text{)} \rightarrow \downarrow \text{error} \rightarrow \text{minimizing mistakes and classifying known data points outweighs a large margin}$

Intercept a_0 : "risk" tolerance, i.e. $\uparrow a_0$ if classifying eating wild mushroom vs not eating it.

$$\text{Range: } \begin{cases} a_0 + m / a_0 - m \\ a_0 + 1 / a_0 - 1 \end{cases}$$

ex) If giving a bad loan is twice as costly as withholding a good loan:

$$\Rightarrow a_1 x_1 + a_2 x_2 + \dots + a_m x_m + \left[\frac{2}{3}(a_0 - 1) + \frac{1}{3}(a_0 + 1)\right] = 0 \rightarrow \frac{1}{3} vs \frac{2}{3} \dots \text{arbitrary}$$

(in soft classification)

Adding a multiplier for each type of error, with a larger penalty the less we want to accept mis-classifying that type of point:

$$\underset{a_0, \dots, a_m}{\text{minimize}} \sum_{j=1}^n m_j \max\left\{0, 1 - \left(\sum_{i=1}^m a_i x_{ij} + a_0\right) y_j\right\} + \lambda \sum_{i=1}^m (a_i)^2$$

$\{ m_j > 1 \text{ for more costly errors}$
 $\{ m_j < 1 \text{ for less costly errors}$

Scaling Data:

We want to ↑margin by $\min_{a_0, \dots, a_m} \sum_{i=1}^m (a_i)^2$, but what if a_m values are magnitudes apart for each attribute?

Example: $\{\text{income in } \$10,000 - \$1,000,000 \text{ range}$ $\Rightarrow \text{income } a_1 (\text{income}) \text{ will dwarf } a_2 (\text{credit score})!$ \Rightarrow solution: Scale data, so that the orders of magnitude are approx the same.

{ Some software automatically do this;
others don't, so do it yourself.

ex) Suppose we try coefficients 5 for income and 701 for credit score, with intercept of a million.

$$\Rightarrow 0 = a_0 + \sum_j a_j x_j \Rightarrow 0 = a_0 + a_1 x_1 + a_2 x_2$$

$$\Rightarrow 0 = 1,000,000 + 5x_{\text{income}} + 701x_{\text{credit}}$$

is the SVM.

Sum of squared coefficients:

$$\Rightarrow \sum_j a_j^2 \Rightarrow 5^2 + 701^2 \Rightarrow 490,025$$

Let's say credit score coefficient increased by one...

$$\Rightarrow 5^2 + 702^2 = 491,426 \rightarrow +101 \text{ increase!}$$

By increase for a small change in coefficient.

... $\dots + 1 \dots + 2 \dots + \dots + n \dots$ is more than one! And a different line!

$$\Rightarrow \sum a_i^2 \Rightarrow 5^2 + 701^2 \Rightarrow 496,025$$

Let's say credit score coefficient increased by one...

$$\Rightarrow 5^2 + 701^2 = 491,426 \rightarrow +101 \text{ increase!}$$

Big increase for a small change in coefficient.

To match this number, x_{score} would have to be 37.8 ... much more than one! And a different line!
(100% change)

Conclusion: SDM model won't work well without scaling.

Adjusting the Data (Scaling)

(Common scaling: data between 0 and 1 ...) → which method to use?

Scaling by factor:

- Let: $\cdot x_{\min j} \rightarrow \min x_{ij}$ (smallest # in data set)
- $\cdot x_{\max j} \rightarrow \max x_{ij}$ (biggest # in data set)

so, for each data point j :

$$x_{ij}^{\text{scaled}} = \frac{x_{ij} - x_{\min j}}{x_{\max j} - x_{\min j}} \xrightarrow{\text{in most}} x_{\text{scaled}} = \frac{x - x_{\min j}}{x_{\max j} - x_{\min j}}$$

IN GENERAL...

i.e. RGB (0-255);
SAT scores (200-800)

Use scaling for bounded ranges,
& standardization for others. → i.e. principal component analysis;
clustering

But sometimes not clear. Try both!

(General) Scaling between b and a :

$$\Rightarrow x_{ij}^{\text{scaled}(a,b)} = x_{ij}^{\text{scaled}(0,1)}(a-1) + b$$

Scaling to a normal distribution ("Standardization")

Common scaling: mean (μ) = 0, standard deviation (σ) = 1

- Factor j has a mean of $\mu_j = \frac{\sum x_{ij}}{n}$
- Factor j has a standard deviation σ_j

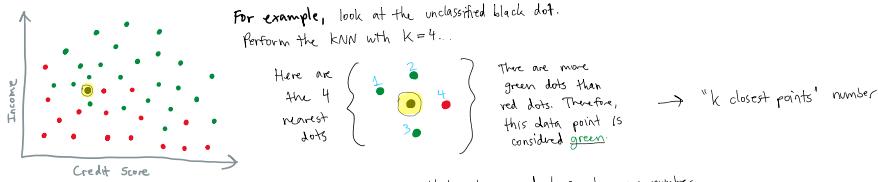
so, for each data point j :

$$x_{ij}^{\text{standardized}} = \frac{x_{ij} - \mu_j}{\sigma_j}$$

Always
Scale/standardize
your data!

k-Nearest Neighbor (kNN) Model

Assume a new data point comes in, and there is no red/green classification on it yet.
Count k number of dots around it, and classify it the same as the greatest dot classification.



Note: can have multiple colors, and k can be any number.

What if it's not clear which dots are closest? → Calculate it!

- Straight-line distance = $\sqrt{\sum_{i=1}^n |x_i - y_i|^2}$

- Weighing attributes by importance: $\sqrt{\sum_{i=1}^n w_i |x_i - y_i|^2}$
weight, set yourself; few for important attributes

What's a good value of k ? ... We'll learn about that! → "Validation" aka measuring the quality of a model



Classification Topic Summary

- Divide data points into groups based on similarity (comes up a lot in analytics)
- Graphical intuition (look at a graph and know classifiers)
- Mathematic Solution Methods
 - Support Vector Machine (SVM)
 - k-nearest-neighbor (kNN)

These are
ML algorithms!

Additional subjects ...

- Data terminology & data types
- Validation
- Distance metrics
- Confusion matrices

Validation — checking how good a model is

- How good is the model?
- How accurate is it?
- How well does it predict X?

Validation does NOT depend on data point accuracy for classifier...

ex) If 24 data points and 3 are incorrectly classified, it is ~~incorrect~~ to assume $21/24 = 87.5\%$ accurate... too optimistic!

Data has two types of patterns:

- 1) Real effect: real relationship between attributes and response } Problem: hard to know which one it is...
- 2) Random effect: random, but looks like a real effect

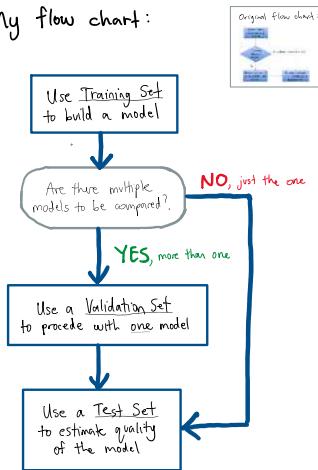
Solution: fit matches on both real & random effects

→ Real effects: same in all data sets

→ Random effects: different in all data sets

Training, Validation, & Test Data Set

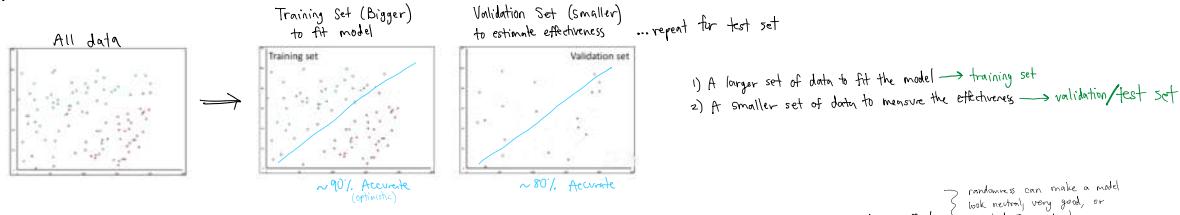
My flow chart:



Need to split data for training, validating, and testing.
Using all data for training is a terrible approach because then there wouldn't be any data left to test the method.
Likewise, reusing same data for training & testing is a bad idea because we need to know how the method will work on data it wasn't trained on.

- Training Set: building models
- Validation Set: picking one model (if applicable)
- Test Set: estimating performance of model

For example.... (let's say multiple models)



Since observed performance = real quality + random effects, a high performing model is more likely to have above-average random effects

so, observed performance of a given model is probably too optimistic.

Splitting: how to split data into training, validation, and test sets

How much data goes in each set? Generally/rule of thumb

- 1 model: 70-90% for training, 10-30% for test → let it equal 100%
- 2+ models: 50-70% for training, and the rest is split equally between validation and test → i.e. 60% training, 20% validation, 20% test

Knowing the total number of data points, how are they split into each set?

ex) Let's say 1000 data points: 60% training, 20% validation, 20% test

Method 1: Random

- Randomly choose 600 data points for training
- Randomly choose 200 (of the remaining 400) data points for validation
- The remaining 200 data points make up the test set

Pros: eliminates bias

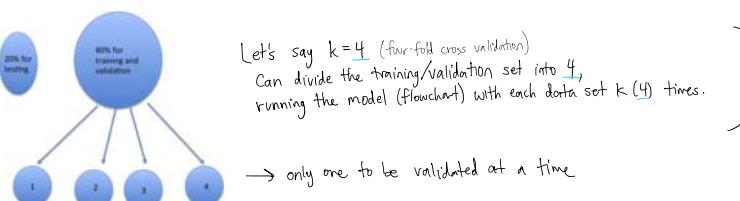
Cons: due to randomness, data can be unevenly separated
i.e. if 10 yrs of data, may have data sets with much more data from specific yearsMethod 2: Rotation

- Take turns selecting points
- i.e. 5 data point rotation sequence could look like:

Training → Validation → Training → Test → Training → repeat

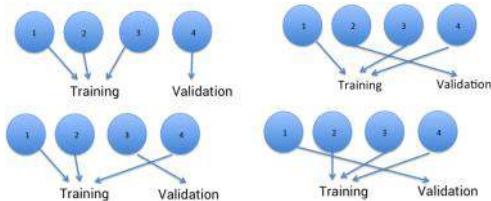
Pros: can ensure data is partitioned evenly
i.e. if 10 yrs of data, can ensure data sets include all yearsCons: can introduce bias
i.e. if 5 pt rotation, and based on Mon-Fri, all Mondays would be in one set, Tuesdays in another, etc.

Cross Validation: Use if important data only appears in the validation or test sets (not training set)
 Denoted as "k-fold cross validation"



No standard number to use as k , but $k=10$ is common.

Visual Example:



Basically, for each of the k parts:

- Train the model on all the other parts
- Evaluate (validate) it on the remaining part

Then, average the k evaluations to estimate the model's quantity.
 Can compare different model types and use this value to pick the best model method!

But for SVM, do NOT average: Instead, train the model again using all the data.

Clustering

→ taking a set of data points and dividing them up into groups; each group contains points that are close to each other or similar.



Examples of use:

- targeted marketing
- personalized medicine
- locating facilities
- image analysis
- Data investigation

Distance Norms

(2-norm) Euclidean (straight-line) distance:

$$\begin{aligned} \Rightarrow \text{distance} &= \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2} \\ \Rightarrow \text{distance} &= \sqrt{|x_1 - y_1|^2 + |x_2 - y_2|^2} \\ &= \sqrt{\sum_{i=1}^n |x_i - y_i|^2} \end{aligned}$$

$\xrightarrow{\text{let } p=1}$

(p-norm) Minkowski distance:

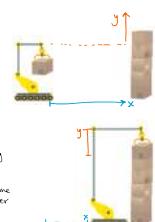
$$\begin{aligned} \Rightarrow \text{distance} &= \sqrt[p]{|x_1 - y_1|^p + |x_2 - y_2|^p} \\ \Rightarrow \text{distance} &= \sqrt[p]{|x_1 - y_1|^p + |x_2 - y_2|^p} \\ &= \sqrt[p]{\sum_{i=1}^n |x_i - y_i|^p} \end{aligned}$$

$\xrightarrow{\text{let } p=\infty}$

∞ -norm distance

$$\begin{aligned} \Rightarrow \text{distance} &= \sqrt[n]{\sum_{i=1}^n |x_i - y_i|^\infty} = \sqrt[n]{\max_{i=1}^n |x_i - y_i|^\infty} \\ &\quad \text{sum equals the largest } |x_i - y_i|^\infty! \\ &= \sqrt[n]{\max_{i=1}^n |x_i - y_i|^\infty} \\ \Rightarrow \max |x_i - y_i| &\quad \checkmark \text{That's it! Just the largest set of numbers in } \rightarrow \text{"the biggest absolute value"} \end{aligned}$$

When would you use this?



k-means algorithm

Clustering

Groupping data points
 Solve using k-means algorithm



x_{ij} = attribute j of data point i

y_{ik} = 1, if data point i is in cluster k
 0, if not

z_{jk} = coordinate j of cluster center k

What we have to find is a set of k cluster centers and assignments of each data point to the cluster it belongs to. This means the total distance from each data point to its cluster center. This plus adds all of the distances from data points to cluster centers, but only when the data point is in that cluster.

$$\text{Minimize } \sum_{i=1}^n \sum_{k=1}^K \sum_j (x_{ij} - z_{jk})^2$$

Subject to $\sum_k y_{ik} = 1$ for each i

Goal: find a set of k cluster centers and coordinates of each data point to a center cluster... all to minimize the total distance of each data point to its center cluster.

This adds the distance from data points to center clusters,

but only when the data point is in that cluster.

Problem! This is a hard optimization problem to solve...

So instead, we use the **k-means algorithm**:

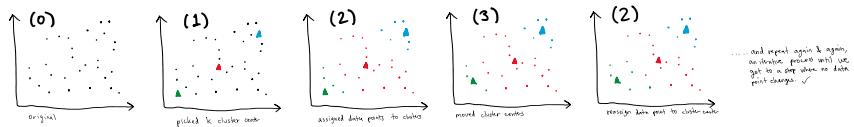
(1) Decide how many clusters we want the algorithm to give us by picking k cluster centers within range of data.

(2) Assign each data point to the nearest cluster center.

(3) But the points we called cluster points aren't really at the center of their clusters, so recalculate cluster centers (centers).

... Repeat, starting with step 2, until no changes

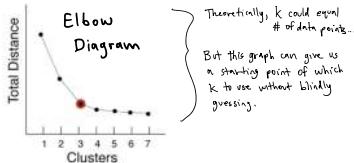
- (1) Decide how many clusters we want to have
- (2) Assign each data point to the nearest cluster center.
- (3) But the points we called cluster points aren't really at the center of their clusters, so recalculate cluster centers. (centroids)
- ... Repeat, starting with step 2, until no changes



Note that the k-means algorithm is:

- a heuristic → it's fast and it works, but it's not guaranteed to find the absolute best solution ... but it comes close!
- an expectation-maximization algorithm ... see steps 1-3 above.
 - Expectation: calculating cluster centers is taking the mean of all the points in the cluster
 - Maximization: reassigning data points to cluster centers = minimizing the negative of the distance to a cluster center
 - minimizing finding the closest distance to a cluster center = maximizing the negative of the distance to a cluster center
 - Takes turns between expectation, maximizing, expectation, maximizing, etc → hence, "em" algorithm

Since k-means algorithm = heuristic = fast, we can run it several times to determine how many clusters centers k to use, comparing the results to determine which to pick.



What to do with outliers ...

Outlier = data point that's not close to any of the cluster centers
The k-means algorithm will assign the outlier to whichever cluster center is closest... but as you can see in visual, this isn't always right.

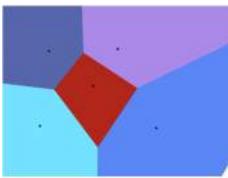
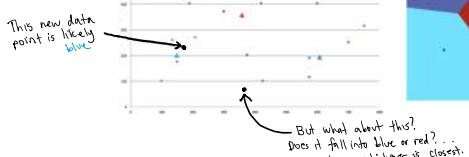
Ask yourself the following:
What is this outlier?
What makes its attribute so different?
What's the implication of putting a point like this in the nearest cluster?

Understand the situation! Look deep at the data, algorithm, & situation to find the right answer.

Predictive Clustering

Based on previous clusters, new data points can be "predicted" into a cluster.

For example ...

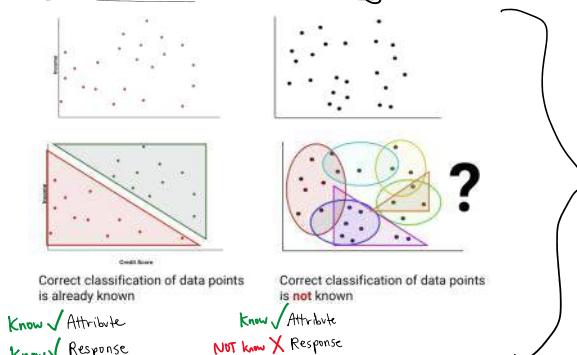


Can color-code the cluster ranges so we know where each data point could fall.

This is called the **Voronoi Diagram**.

But what about this?
Does it fall into blue or red?
→ choose whichever is closest.

Classification vs. Clustering



Supervised Learning: (classification)

machine learning where the "correct" answer/response is known for each data point in the training set → more common

Unsupervised Learning: (clustering)

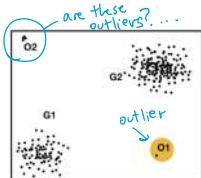
machine learning where the "correct" answer/response is NOT known for each data point in the training set

Week 3: Basic Data Preparation, Change Detection

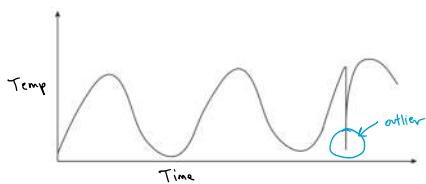
Tuesday, February 1, 2022 4:02 PM

How to Prepare a set of data for use in analytics.

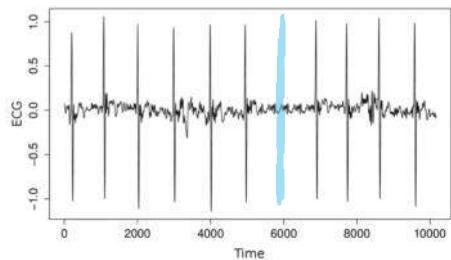
Data Preparation → process data before building model... aka address problems in data beforehand
i.e. scale data, look for outliers, etc.



Outlier: data point that's very different from the rest

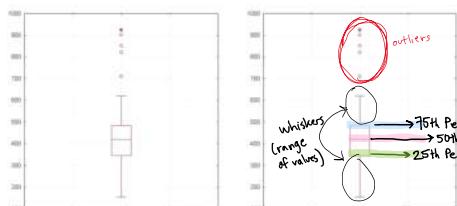


Contextual Outlier: $f(t) = y$ value isn't an outlier since it's not far from other points overall, but the value at time t is an outlier since it's far from points nearby in time



Collective Outlier: outlier by omission; something is missing in a range of points, but can't tell exactly where

Noticing outliers is easy. It's figuring out what to do with them that's hard.



Box-and-whisker Plot

... to find outliers in one dimension

{ For example ...

Say we did exponential smoothing here.
At that outlier point, the model's error will be very high.

Unfortunately, there's no great way to detect multi-dimensional outliers.

However!

We can build a model, fit the parameters, then see which points have ↑ error.

Outliers: real, fake, or random?

Sometimes, the outliers are simply bad data. Other times, something is wrong with the situation/data set.
But most of all, it could just be randomness. 4% of data outside two standard deviations are expected to be outliers.

Consider context to outliers. There is no one right answer.

Change Detection → determining whether something has changed (usually time series data)

Why is this needed? →

- determine if action might be needed i.e. time for preventive maintenance? global temperature change? decreased sales?
- determine impact of past action i.e. policy change have desired effect? with lower tax rates, were more people hired? did discount affect sales?

- Why is this needed? →
- determine if action might be needed i.e. time for preventive maintenance? global temperature change? decreased sales?
 - determine impact of past action i.e. policy change have desired effect? with lower tax rates, were more people hired? did discount affect sales?
 - determine changes to help plan i.e. have voting patterns changed?

Case Studies

Manufacturing plant: can predict when there is about to be a bottleneck, to prevent bottlenecks

Railroad Safety: to prevent train derailment, temperature sensors are along the track to predict failure (derailment). If it's too high, the train will stop before derailment, do maintenance, then resume.

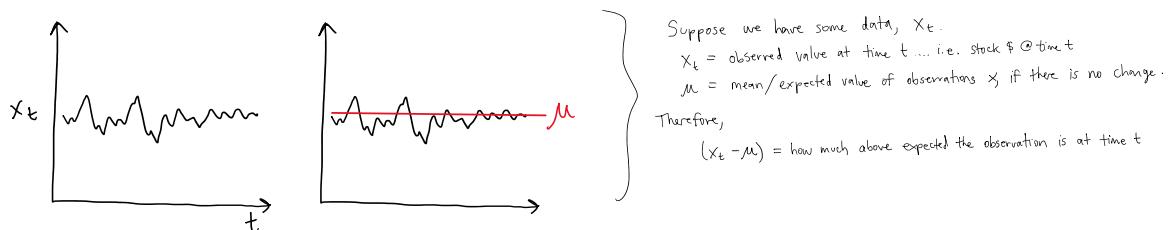
Problem: sometimes there are outliers... we wouldn't want to stop the train if it's not going to fail.
(random fluctuations)

So, how do we determine if there's actually a change? Need something faster than hypothesis testing....
That's where CUSUM comes in.

Cumulative Sum (CUSUM)

CUSUM detects increases, decreases, or both.

Answers the question: has the mean of the observed distribution gone beyond a critical level?



CUSUM Metric: Calculate S_t and declare we've observed a change if S_t goes above threshold T .

$$S_t = \max \{ 0, S_{t-1} + (x_t - \mu) \}$$

Is $S_t \geq T$?
we pick threshold T

Choose zero if $S_{t-1} + (x_t - \mu)$ is negative
adding to the previous value
amount above μ baseline line

HOWEVER!

Basically, if the running total would be below zero, it's irrelevant to the question of whether we later see an increase

Sometimes there is randomness.
Sometimes, x_t will randomly be higher than the expectation.

How do we account for this?

$$S_t = \max \{ 0, S_{t-1} + (x_t - \mu - c) \}$$

Is $S_t \geq T$?
we pick threshold T

Choose zero if $S_{t-1} + (x_t - \mu - c)$ is negative
adding to the previous value
amount above μ baseline line
sensitivity value to account for randomness (like standard deviation)

$\uparrow c, \downarrow S_t, \uparrow \text{sensitive}$

or

$\downarrow c, \uparrow S_t, \downarrow \text{sensitive}$

$\uparrow T, \downarrow \text{time to detect changes}, \downarrow \text{false changes}$

or

$\uparrow T, \uparrow \text{time to detect changes}, \uparrow \text{false changes}$

Detecting an increase

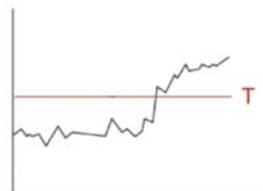
$$S_t = \max \{ 0, S_{t-1} + (x_t - \mu - c) \}$$

Is $S_t \geq T$?

Detecting a decrease

$$S_t = \max \{ 0, S_{t-1} + (x_t + \mu - c) \}$$

Is $S_t \geq T$?



model parameters

$\dots \overbrace{- \quad - \quad - \quad -}^{\text{model parameters}}$

model parameters
How do you select C & T ?

Model needs to find changes quickly, but not be so sensitive that it says there's a change when there isn't one.
But the parameters are totally dependent on risk tolerance of the situation.

...
play around with T until it gets expected results?
 C can be 2-3 standard deviations, but it still totally depends.

Time Series Data: data where the same response is known for many time periods

Still, there could be trends over time, cyclical variations, and randomness

Say there's a patient getting a blood pressure reading, and it comes out higher than average

... But this could be random fluctuation. Should the doctor take action?

without knowing a baseline value, it's hard to know what to do. This is where exponential smoothing comes in handy.

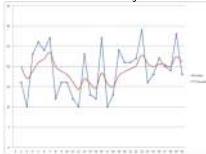
Let's go back to the blood pressure scenario.

Say:

S_t = expected baseline response at time period t → blood pressure at hour t

x_t = the observed response → observed blood pressure at hour t

Exponential Smoothing: A visual



2 possibilities. Either:

$S_t = x_t \rightarrow$ the observed value (blood pressure) is the real indicator of the baseline as in, blood pressure is always that high.

$S_t = S_{t-1} \rightarrow$ today's baseline is the same as yesterday's baseline as in, blood pressure is unusually high.

Exponential Smoothing Method combines these two ideas.

Other names:
Single/Double/Triple Exponential Smoothing
Holt-Winters Method
Smoothing

$$S_t = \alpha x_t + (1-\alpha)S_{t-1}$$

[single]

Initial condition: $S_1 = x_1$

α is a number from 0 to 1

We choose α :
 $0 < \alpha < 1$
 $\Rightarrow \alpha \rightarrow 0$: a lot of randomness → trusting S_{t-1} (previous estimate)
 $\Rightarrow \alpha \rightarrow 1$: little randomness → trusting x_t (observed value/what you see)

Note that every new observation contributes to the current baseline estimate.
It's just baked in the S_{t+1} equation.
But, newer observations are more important/has more weight.

What if there's a trend? The value we're measuring may naturally be increasing or decreasing over time.

Use this equation instead:

$$S_t = \alpha x_t + (1-\alpha)(S_{t-1} + T_{t-1})$$

$T_t = (\beta)(S_t - S_{t-1}) + (1-\beta)(T_{t-1})$

α is a number from 0 to 1
 β is a number from 0 to 1

[trend/double]

$$S_t = \alpha x_t + (1-\alpha)(S_{t-1} + T_{t-1})$$

$$T_t = (\beta)(S_t - S_{t-1}) + (1-\beta)(T_{t-1})$$

Initial condition: $T_1 = 0$

We choose α, β :
 $0 < \alpha < 1$
 $0 < \beta < 1$

What if there's a cyclic pattern? Trends may go up and down, but they repeat over a cycle/set length of time L .

You could add up trends in an additive way, or use equation below.
cyclic pattern = seasonalities in a multiplicative way

Use this equation instead:

$$C_t = \gamma \left(\frac{x_t}{S_t} \right) + (1-\gamma)C_{t-L}$$

$$S_t = \frac{\alpha x_t}{C_{t-L}} + (1-\alpha)(S_{t-1} + T_{t-1})$$

γ is a number from 0 to 1
 α is a number from 0 to 1
 L is the length of cycle ... i.e. same day next year or in a day
 γ is the cyclic factor ... i.e. same day next year or in a day

[cyclical/multiplicative seasonality/triple]

$$S_t = \frac{\alpha x_t}{C_{t-L}} + (1-\alpha)(S_{t-1} + T_{t-1})$$

$$C_t = \gamma \left(\frac{x_t}{S_t} \right) + (1-\gamma)C_{t-L}$$

Initial condition: first L values of C are set to 1
... since $C=1$ equals no cycle, easier to start

Forecasting

Exponential smoothing can analyze what happened in the past, but it can also predict what will happen next.

Basic exponential smoothing equation: $S_t = \alpha x_t + (1-\alpha)S_{t-1}$

Prediction: $S_{t+1} = \alpha x_{t+1} + (1-\alpha)S_{t+1-1}$

x_{t+1} is unknown ...
so, best guess: $x_{t+1} = S_t \Rightarrow$ "I predict the next observed value will be the growing baseline estimate."

S_t (instead of x_{t+1})

Best Guess
 $x_{t+1} = S_t$

Therefore, forecast for time period $t+1$:

$$S_{t+1} = \alpha x_{t+1} + (1-\alpha)S_t$$

$$S_{t+1} = F_{t+1} = \alpha S_t + (1-\alpha)S_t$$

Let $\alpha = 0$...

$$\hookrightarrow F_{t+1} = \gamma(\alpha)S_t + (1-\gamma)\alpha S_t$$

[single - forecasting]

$$F_{t+1} = S_t$$

$$F_{t+k} = S_t, k=1,2,3\dots$$

Trend Exponential Smoothing Equation:

$$S_t \Rightarrow S_t = \alpha x_t + (1-\alpha)(S_{t-1} + T_{t-1})$$

$$S_{t+1} \Rightarrow S_{t+1} = \alpha x_{t+1} + (1-\alpha)(S_t + T_t)$$

γ is the alpha again

$$\gamma = \alpha S_t + (1-\alpha)(S_t + T_t)$$

[double/trends - forecasting]

$$F_{t+1} = S_t + T_t$$

Best Guess
 γ

Best Guess

$$S_t \Rightarrow S_t = \alpha X_t + (1-\alpha)(S_{t-1} + T_{t-1})$$

$$\text{All } \alpha \rightarrow 1 \Rightarrow S_{t+1} = \alpha X_{t+1} + (1-\alpha)(S_t + T_t)$$

↑ is the exp. Sm. again

$$\begin{aligned} \text{Simplifying, } &= \alpha S_t + (1-\alpha)(S_t + T_t) \\ \text{with } X_t = S_t + T_t &= \alpha S_t + S_t + T_t - \alpha S_t - \alpha T_t \\ &= S_t + T_t - \alpha T_t \end{aligned}$$

..... Let $\alpha = 0 \dots$

$$S_{t+1} = F_{t+1} = S_t + T_t$$

[double/trends - forecasting]

$$F_{t+1} = S_t + T_t$$

$$F_{t+k} = S_t + kT_t, k=1,2,3\dots$$

Multiplicative Seasonality Exponential Smoothing Equation:

$$C_t = \gamma \left(\frac{X_t}{S_t} \right) + (1-\gamma) C_{t-1}$$

$$C_{t+1} = \gamma \left(\frac{X_{t+1}}{S_{t+1}} \right) + (1-\gamma) C_{t+1-L}$$

..... Let $\gamma = 0 \dots$

$$\Rightarrow C_{t+1} = C_{t+1-L}$$

$$\Rightarrow S_t = \frac{\alpha X_t}{C_{t-L}} + (1-\alpha)(S_{t-1} + T_{t-1})$$

$$\Rightarrow S_{t+1} = \frac{\alpha X_{t+1}}{C_{t+1-L}} + (1-\alpha)(S_t + T_t)$$

$$= \frac{\alpha S_t}{C_{t+1-L}} + (1-\alpha)(S_t + T_t)$$

..... Let $\alpha = 0 \dots$

$$\Rightarrow S_{t+1} = F_{t+1} = \frac{(0)S_t}{C_{t+1-L}} + (1-\alpha)(S_t + T_t)$$

$$F_{t+1} = S_t + T_t \quad ???$$

[triple/cyclical - forecasting]

$$F_{t+1} = (S_t + T_t) C_{(t+1)-L}$$

$$F_{t+k} = (S_t + kT_t) C_{(t+1)-L-(k-1)}, k=1,2,3\dots$$

Which α , β , and γ do I use?

Use Optimization! Use sum of squared error equation: $(F_t - X_t)^2$

Use the α, β, γ values in which error is minimized.



Exponential smoothing is for analyzing time series data where the same response is known for many time periods.

What if I want something more general?

ARIMA : AutoRegressive Integrated Moving Average

↪ Estimates or forecasts a value.

- A more general way to analyze time series data.
- Short-term forecasting...when data is more stable with fewer peaks, valleys, & outliers
- Need 40 past data points for ARIMA to work well

3 parts: (1) Differences, (2) Autoregression, (3) Moving Parts

(1) Differences

We know that the exponential smoothing equation is:

$$S_t = \alpha X_t + (1-\alpha) S_{t-1}$$

$$\Rightarrow S_t = \alpha X_t + (1-\alpha) \alpha X_{t-1} + (1-\alpha)^2 \alpha X_{t-2} + (1-\alpha)^3 \alpha X_{t-3} + \dots + (1-\alpha)^n \alpha X_{t-n}$$

This works well if the data is stationary.
meaning that the mean, variance, and other measures are all expected to be constant over time.
..... But what if it's not stationary? i.e. has trend or seasonality

Often, data is not stationary... but the differences are.

First-order difference $D_{(1)}:$ difference of consecutive observations
 $D_{(1)t} = (x_t - x_{t-1})$
 Second-order difference $D_{(2)}:$ differences of the differences
 $D_{(2)t} = (x_t - x_{t-1}) - (x_{t-1} - x_{t-2})$
 Third-order difference $D_{(3)}:$ differences of the differences of the differences
 $D_{(3)t} = [(x_t - x_{t-1}) - (x_{t-1} - x_{t-2})] - [(x_{t-1} - x_{t-2}) - (x_{t-2} - x_{t-3})]$
 ↳ dth-order differences $\rightarrow D_{(d)}$

(2) Autoregression

Exponential Smoothing is an order- ∞ autoregressive model.
It uses data all the way back.

→ Order-p autoregressive models only go back p time periods.

⇒ ARIMA combines autoregression and differencing → autoregression on the differences

Use p time periods of previous observations
to predict d-th order time differences.

(3) Moving Average

Previous errors E_t as predictors: $E_t = \hat{x}_t - x_t$

If we go back q time periods, it's called an order-q moving average.

ARIMA(p,d,q) Model

$$D_{(d)t} = \mu + \sum_{i=1}^p \alpha_i D_{(d)t-i} - \sum_{i=1}^q \theta_i (\hat{x}_{t-i} - x_{t-i})$$

dth order differences
 pth-order autoregression
 qth-order moving average

Statistical software can be used to find d,p,q.

But also, here are some common ARIMA(p,d,q) functions:

- ARIMA (0,0,0): white noise (no pattern)
- ARIMA (0,1,0): "random walk"
- ARIMA (p,0,0): AR (Autoregressive) model
- ARIMA (0,0,q): MA (Moving Average) model
- ARIMA (0,1,1): basic exponential smoothing model

ARIMA is great for estimating or forecasting a value.

But what if you want to estimate or forecast a variance? → use GARCH

GARCH: Generalized Autoregressive Conditional Heteroskedasticity

↳ Estimate or forecasts the variance.

Why would we need to know the variance?

⇒ variance = how much forecast might be higher or lower than the true value

So, ⇒ Using GARCH can help us estimate the amount of error

Variance estimation is important in investment portfolio optimization model
 risk vs. returns

Variance: a proxy for the amount of volatility or risk
 risk vs. returns

GARCH (p,q) Model

$$\sigma_t^2 = \omega + \sum_{i=1}^p \beta_i \sigma_{t-i}^2 + \sum_{i=1}^q \gamma_i \epsilon_{t-i}^2$$

Similar to ARIMA, but with key differences:

ARIMA	GARCH
• Observations	• Variances
• Linear Errors	• Squared Errors
• Differences of Variances	• Raw Variances
• Has d value: ARIMA(p,d,t)	• No d value: GARCH(p,q)

[summary]



GARCH

$$\sigma_t^2 = \omega + \sum_{i=1}^p \beta_i \sigma_{t-i}^2 + \sum_{i=1}^q \gamma_i \epsilon_{t-i}^2$$

ARIMA

$$D_{(d)t} = \mu + \sum_{i=1}^p \alpha_i D_{(d)t-i} - \sum_{i=1}^q \theta_i (\hat{x}_{t-i} - x_{t-i})$$

Week 5: Basic Regression

Friday, March 4, 2022 3:42 PM

Regression is a way to determine if a variable has an impact on the system.

Regression can:

- (1) Answer questions about how systems work (descriptive)
- (2) Make predictions about what will happen in the future (predictive)

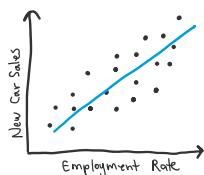
- EXAMPLES
- (1) how many runs the average home run is worth in baseball
 - (2) estimate the effect of economic factors on presidential election voting
 - (3) that the influence of a general education on other future knowns
 - (4) discover key factors in purchasing a specific type of car
 - (5) predict how fast a child will be as an adult
 - (6) forecast the price of oil a year into the future
 - (7) predict the remaining lifetime of someone applying for life insurance
 - (8) predict demand & pricing in a city over the next six months

Simple Linear Regression (SLR)

Linear Regression with one predictor.

Looking for a linear relationship between the predictors and the response.

Let's look at an example....



$$y = \text{response} \rightarrow \text{new car sales}$$

$$x_1 = \text{predictor} \rightarrow \text{employment rate}$$

So, when more people are working, there are more people who can afford to buy a new car, so new car sales are higher.... But how much higher?? That's what regression will find out.

Regression Equation: $y = a_0 + a_1 x_1$

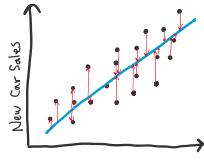
But what if we had m predictors?....

Regression Equation:

$$\begin{aligned} &\text{with one predictor: } y = a_0 + a_1 x_1 \\ &\text{with } m \text{ predictors: } y = a_0 + a_1 x_1 + a_2 x_2 + \dots + a_m x_m \\ &\Rightarrow y = a_0 + \sum_{j=1}^m a_j x_j \end{aligned}$$

Is this a good line?....

Find the sum of squared errors to measure the quality of the line's fit.



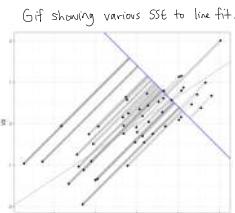
$$y_i = \text{cars sold at data point } i$$

$$\hat{y}_i = \text{model's prediction of cars sold (@ regression line of } f(x_i))$$

$$\begin{aligned} \text{Data point } i \text{ prediction error: } y_i - \hat{y}_i &= y_i - (a_0 + a_1 x_{ii}) \\ &\Rightarrow \hat{y}_i = a_0 + a_1 x_{ii} \end{aligned}$$

Employment Rate

Sum of Squared Errors: $\sum_{i=1}^n (y_i - \hat{y}_i)^2$



Minimize SSE to find the best Regression Line

$$\begin{aligned} \text{SSE} &= \sum_{i=1}^n [y_i - (a_0 + a_1 x_{ii})]^2 \\ &= \sum_{i=1}^n [y_i - (a_0 + \sum_{j=1}^m a_j x_{ij})]^2 \end{aligned}$$

Underlying math: SSE is a convex function we're trying to minimize. Take partial derivatives of the SSE wrt each constant. Set each partial derivative to zero, then solve those equations simultaneously to find the minimized SSE, aka best fit solution → or just use software.

How do we measure the quality of a model's fit?

A few different ways. Examples: Likelihood, AIC, BIC

Likelihood: measures the probability (density) for any parameter set

We assume that the observed data is the correct value, and that we have information about the variance.

Maximum Likelihood: parameters that give the highest probability → aka, best fit set of parameters

ex) Error $\sim N(0, \sigma^2)$, i.i.d.
Observations: z_1, \dots, z_n
Model estimates: y_1, \dots, y_n

Suppose we believe that errors are normally distributed with mean zero and variance σ^2 , and they're independent from one data point to the next.

Then, given observed data z_1 through z_n , and our model estimates y_1 through y_n , below is the probability of observing z_i if the true value is y_i .

ex) Error $\sim N(0, \sigma^2)$
 Observations: z_1, \dots, z_n
 Model estimates: y_1, \dots, y_n

normally distributed with mean zero and variance σ^2 , and they're independent from one data point to the next.

estimates y_i through y_n , below is the probability of observing z_i if the true value is y_i .

Probability density for observing z_i if true value is y_i

$$= \left(\frac{1}{\sigma \sqrt{2\pi}} \right) \left(e^{-\frac{(z_i - y_i)^2}{2\sigma^2}} \right)$$

since independent?

$$\Rightarrow \left(\frac{1}{\sigma \sqrt{2\pi}} \right)^n \left[e^{-\frac{1}{2\sigma^2} (\sum_{i=1}^n (z_i - y_i)^2)} \right]$$

Can find largest value by finding n

Simplifying this gets us MLE
 \Rightarrow Maximum Likelihood Estimation

Maximum Likelihood Estimation (MLE):

the set of parameters that minimizes the sum of squared errors

$$\text{maximize } \left(\frac{1}{\sigma \sqrt{2\pi}} \right)^n e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (z_i - y_i)^2} \rightarrow \left[\text{minimize } \sum_{i=1}^n (z_i - y_i)^2 \right]$$

Therefore, in regression:

$$SSE = \sum_{i=1}^n [y_i - (a_0 + \sum_{j=1}^m a_j x_{ij})]^2$$

Remember, this assumes model is independent with normally distributed errors

$$\boxed{\text{MLE} = \text{minimize} \sum_{i=1}^n [z_i - (a_0 + \sum_{j=1}^m a_j x_{ij})]^2}$$

... y_i & z_i are the same b/w

Note that we can compare two different models by using the likelihood ratio (ratio of their likelihood), then conduct a hypothesis test.

Some more complex approaches...

↑ number of parameters in model ... we also want to prevent overfitting. So either do well to fit random factors or the real ones.

Akaike Information Criterion (AIC):

$$\Rightarrow AIC = 2k - 2 \ln(L^*)$$

"Penalty Term"
 Balances likelihood with simplicity to help prevent overfitting.

L^* = maximum likelihood value
 k = number of parameters being tested (attributes)

$$AIC = 2k - 2 \ln(L^*)$$

Now to substitute this into regression.

$$AIC = 2k - 2 \ln(L^*)$$

$m+1$, for a_0 there are m others

Sub in likelihood function

$$\Rightarrow AIC = 2(m+1) - 2 \ln \left(\left(\frac{1}{\sigma \sqrt{2\pi}} \right)^n e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (z_i - (a_0 + \sum_{j=1}^m a_j x_{ij}))^2} \right)$$

↓ AIC is preferred;
 results in ↓k and ↑ L^*

Note that AIC has nice properties, if there are infinitely many data points.

Obviously this is impossible ... so, for "smaller" data sets, use Corrected AIC (AIC_c)

Corrected AIC (AIC_c): we generally use this unless we have an extremely large ($\rightarrow \infty$) data set

$$\Rightarrow AIC_c = AIC + \frac{2k(k+1)}{n-k-1} \Rightarrow AIC_c = 2k - 2 \ln(L^*) + \frac{2k(k+1)}{n-k-1}$$

When comparing models with different AIC_c /AIC values, we can calculate the relative probability that one of them is better than the other.

When comparing models with different AIC_c/AIC values, we can calculate the relative probability that one of them is better than the other.

This is called Relative Likelihood.

$$\text{Relative Likelihood} = e^{\frac{AIC_1 - AIC_2}{2}}$$

Example:
 Model 1: AIC = 75
 Model 2: AIC = 80
 $\Rightarrow \text{Relative Likelihood} = e^{\frac{(75-80)}{2}} = 8.2\%$
 ... Model 2 is 8.2% as likely as Model 1 to be better... that's pretty low.
 So, Model 1 is probably better and thus preferred.
 ↳ we know this already since AIC_c is smaller...?

But what if you have a lot more data than parameters?

Use Bayesian Information Criterion (BIC).

Bayesian Information Criterion (BIC): similar to AIC, but encourages models with fewer parameters (attributes)

$$\begin{aligned} AIC &= 2k - 2\ln(L^*) \\ \Rightarrow BIC &= k \ln(n) - 2\ln(L^*) \end{aligned}$$

L^* = maximum likelihood value
 k = number of parameters being tested
 n = number of data points

Note that BIC's penalty term > AIC's penalty term...

so, only use BIC when there are more data points than parameters.

$$BIC = k \ln(n) - 2\ln(L^*)$$

When comparing Model 1 BIC (BIC₁) and Model 2 BIC (BIC₂)

- $0 < |BIC_1 - BIC_2| < 2 \Rightarrow$ Smaller-BIC model is "slightly likely" better
- $2 < |BIC_1 - BIC_2| < 6 \Rightarrow$ Smaller-BIC model is "somewhat likely" better
- $6 < |BIC_1 - BIC_2| < 10 \Rightarrow$ Smaller-BIC model is "likely" better
- $|BIC_1 - BIC_2| > 10 \Rightarrow$ Smaller-BIC model is "very likely" better

Final Notes: There is no rule to pick AIC over BIC or MLE.
 All three can give valuable information,
 and using all three can help you decide which model is better.

↑ frequentist POV
 ↓ Bayesian POV

Interpreting Regression Coefficients

Components of Analytics:

- Descriptive Analytics
 - Predictive Analytics
 - Prescriptive Analytics
- Regression can do these two, hence why it's so common.

If we're using descriptive regression to figure out how the system works, the key is to look at the coefficients:

ex) Baseball example: want to determine average number of runs a homerun is worth

Response: how many runs a team scored that season

Predictors: Number of home runs, triples, doubles, singles, outs, double plays, stolen bases, etc.

Regression Model:

$$\text{Runs Scored} = a_0 + a_1[\text{number of home runs}] + a_2[\text{number of triples}] + \dots + a_7[\text{number of stolen bases}]$$

Let's say this is 1.4.
 so, every home run hit will add 1.4 runs to the team's total.

If we're using predictive regression to make forecasts, the key is to look at the predicted responses.

ex) Height Example: want to predict how tall a 2-year-old will be as an adult

Response: a person's height as an adult

Predictors: father's height, mother's height, age 2 height, male/female

Regression Model:

$$\text{Adult Height} = a_0 + a_1[\text{father's height}] + a_2[\text{mother's height}] + \dots + a_4[\text{male or female}]$$

Let's say this is 1.4.
So, every cm of the father's height will add 1.4cm to the predicted adult height

Causation vs. Correlation

Causation: one thing causes another thing

Correlation: two things tend to happen or not happen together

Neither of them might cause the other.

Correlation does NOT equal causation!

ex)

$$y : \text{hours/day spent outdoors in winter}$$

$$x_1 : \text{city's avg. daily winter temperature}$$

$$y = a_0 + a_1 x_1$$

$$\text{Correlation between } y \text{ and } x_1$$

- Low p-value of a_1

Does higher winter temperature cause people to go outside?

Kind of (probably)

$$y : \text{hours/day spent outdoors in winter}$$

$$x_1 : \text{city's avg daily winter temperature}$$

$$x_1 = b_0 + b_1 y$$

$$\text{Same correlation between } y \text{ and } x_1$$

- p-value of b_1 = p-value of a_1

Does people spending more time outdoors in the winter cause higher winter temperatures?

That's pretty silly!

There's a correlation here between temperature and people who go outside, but we can't say one caused the other.

i.e.

X ↑ temperatures ⇒ ↑ people who go outside ... seems possible

X ↑ people who go outside ⇒ ↑ temperatures ... obviously wrong

How to tell when there is causation...

Actually very difficult; easier to prove that there is no causation.

Here are starting points, though:

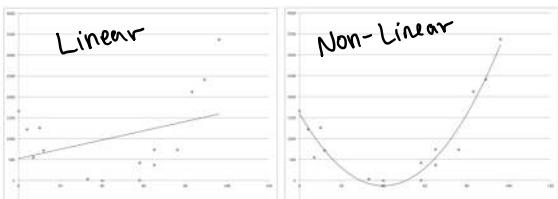
- Cause is before the effect
- The idea of causation makes sense
- No outside factors causing the relationship → most important one, but the hardest. There could be infinite of factors...

In short:
be careful before causing causation.

Generalized Regression Modeling

SLR is great for linear models....

But what if it isn't linear?



Transform the data.

We can adjust the data so the fit is linear.

Quadratic Regression: $y = a_0 + a_1 x_1 + a_2 x_1^2$ → add exponential

Variable Interaction: $y = a_0 + a_1 x_1 + a_2 x_2 + a_3(x_1 x_2)$

Response Transform: $\log(y) = a_0 + a_1 x_1 + \dots + a_m x_m$

Can really transform anything....

$$\text{ex)} \quad y = a_0 + a_1 x_1^{1.5} + a_2 x_1^2 x_3 + a_3 \sin(x_2)$$

$$\ln(y) = a_0 + a_1 x_1^2 x_3 + a_2 \log(x_3)$$



Outputs

So, now we know how to fit and build regression models.
Now it's time to look at outputs and determine which attributes are important, and which are not.

P-Values..... estimates the probability that the coefficient a_m is equal to zero.

(form of hypothesis testing)

i.e. p-value = 0.02 \Rightarrow 2% of not being significant
 \Rightarrow 98% of being significant

Correlations : $\begin{cases} \uparrow p\text{-value}_{(\text{toward } 1)} & \Rightarrow \text{attribute is not important} \\ \downarrow p\text{-value}_{(\text{toward } 0)} & \Rightarrow \text{attribute is important} \end{cases}$

\Rightarrow If $p\text{-value} > 0.05$, then remove the attribute from the model.

Note that other thresholds α can be used.
 $\alpha = 0.05$ is pretty standard though.
 $\cdot \uparrow \alpha, \uparrow \text{attributes}, \uparrow \text{likelihood of irrelevant attribute}$
 $\cdot \downarrow \alpha, \downarrow \text{attributes}, \downarrow \text{likelihood of irrelevant attribute}$

Warning! As $\uparrow \text{data}$, $\downarrow p\text{-values}$... even with irrelevant attributes.
Be careful to rely on p-values when there is a lot of data.

Confidence Interval

Most software will give a 95% confidence interval around that coefficient.
So, you can see where the coefficient probably lies, and how close that is to zero.

t - Statistic

..... the coefficient divided by its standard error

Similar to p-value. It's just another way of finding the same information.

Coefficient

Sometimes you'll discover that the coefficient doesn't make much of a difference,
even when it has a low p-value, and even when multiplied by the attribute.
i.e. when looking at population predictions, an $a_n = 1$ will hardly make a difference.

R-squared Value (R^2)

..... estimates how much variability your model accounts for.

For example, an $R^2 = 59\%$ means the model has accounted for 59% of the variability in the data.
The 41% is from randomness or other factors not modeled.

Adjusted R-squared (R_{adj}^2) : adjusts for the number of attributes used

When using R^2 , keep in mind that some things are not easily modeled.

In homework, we may get R^2 values at the ~90% range.

This hardly reflects real life.

In social sciences, for example, R^2 values of 0.3 or 0.4 are quite good.



Week 6: Advanced Data Preparation

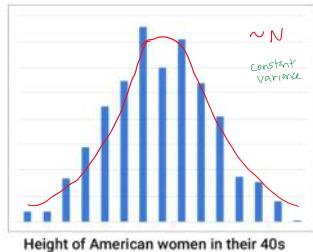
Saturday, March 5, 2022 12:47 PM

Sometimes, it's helpful to transform data before using it to fit a model.

For example....

Recall variance:

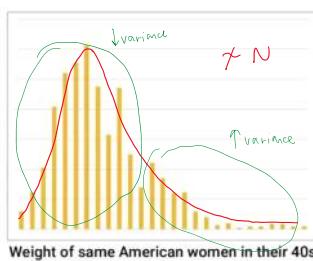
- $\downarrow \sigma^2 \Rightarrow$ data points are similar; do not vary widely from the mean μ
- $\uparrow \sigma^2 \Rightarrow$ data points are not similar; varies widely from the mean μ



Let's say we have a normality assumption.

This would be true for the graph on the left.

Normality is a pretty good fit to the normal distribution, so the variance is constant across the whole distribution.



But not all data has normality.

Results will have bias when normality assumption is wrong.

The graph on the left has varying amount of variance, which can introduce bias.

Heteroscedasticity: unequal / differences in variance

So, what can we do?

How do we deal with heteroscedasticity?

Transform the data to fit the model, i.e. fit to the normal distribution

Box-Cox Transformation transforms response to eliminate heteroscedasticity.

↳ A logarithmic transformation that:

- stretches out the smaller range to enlarge its variability,
- shrinks the larger range to reduce its variability

$$t(y) = \frac{(y^\lambda - 1)}{\lambda}$$

$\lambda \Rightarrow$ value we're optimizing
 $y =$ vector of responses
 $t(y) =$ transform vector

The above equation will result in being close to normal distribution.

When would we have to use the Box-Cox transformation?

↳ Check by using a Q-Q plot.

See Week 9 for more information

Detrending

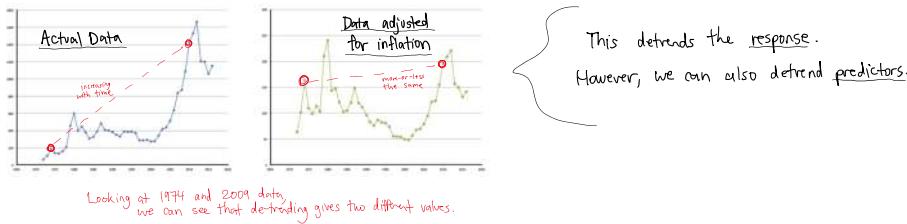
Sometimes trends in time series can influence/mess up a factor-based analysis.

So, we'd want to "detrend"/remove the trend from the data.

Consider detrending when doing factor-based modeling like regression and SVM.

For example...

The Price of Gold over the past 45 years



How to Detrend

In the case of the price of gold example, detrending is easy since there are historical inflation rates.
So, adjust each data point individually by the amount of historical inflation.

However, most cases aren't like this.

In most cases, we don't know what the trend has been in the past, except for our own data we're trying to detrend.

Solution → go factor-by-factor and fit
a one-dimensional regression to it (could use more complicated function if you want)

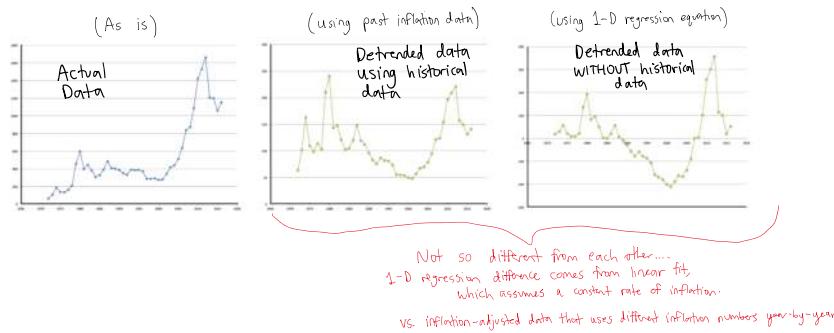
$$\text{One-dimensional regression: } y = a_0 + a_1 x$$

ex) Simple linear regression for gold prices:

$$\begin{aligned} \Rightarrow y &= a_0 + a_1 x \\ \Rightarrow \text{Price} &= -45,600 + (23.2)(\text{year}) \end{aligned}$$

Detrended price:

$$\begin{aligned} \Rightarrow y_{\text{detrend}} &= a_{\text{actual value}} - y \\ \Rightarrow \text{Price}_{\text{detrend}} &= \text{Actual price} - (-45,600 + (23.2)(\text{year})) \end{aligned}$$



Principal Component Analysis (PCA)

Used in high-dimensional and correlated data.

→ PCA removes these correlations and ranks coordinates by importance.

In other words,

A PCA plot converts the correlations (or lack thereof) of all the attributes into a 2D graph.
Cells that are highly correlated are clustered together.

ex) Does the past performance of any particular stock indicate how well the overall market will perform the next day?

→ Say there are 6000+ securities, and we have "normal" stock trade data.
The table could look like this:

↳ Say there are 6000+ securities, and we have "normal" stock trade data.

The table could look like this:

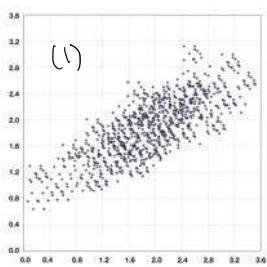
	Stock ₁	Stock ₂	Stock ₃	Stock ₄	Stock ₅	...	Stock ₆₀₀₀
t ₁	~	~	~	~	~	...	~
t ₂	~	~	~	~	~	...	~
t ₃	~	~	~	~	~	...	~
t ₄	~	~	~	~	~	...	~
t ₅	~	~	~	~	~	...	~
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
t _n	~	~	~	~	~	...	~

We can compare Stock₁ and Stock₂, or really Stock₁ with any of them. And same with any two stocks. But how do we compare three stocks? Or four, or five, or six? Can't make 6-D graphs, or a million 2D graphs... This is why PCA is useful.

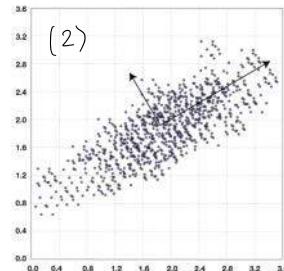
Plus, underlying situations may have changed in time. And there are high correlations between some of the predictors.

→ PCA can "reduce" the amount of data needed by addressing both of these issues.

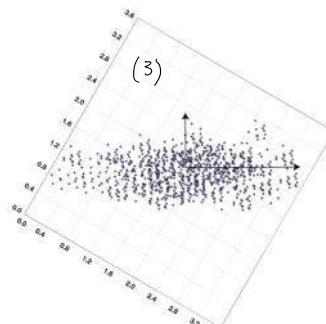
What does "removing correlation" look like? A visual



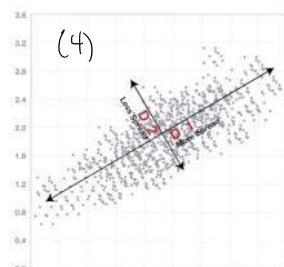
(1) There is clear correlation here...



(2) Marking the correlation



(3) If we pivot the graph, now the marked correlation line is at (0,0). So, the correlation is removed!
This is just to visually see what's going on.
We can rotate it back and redefine a coordinate system.



(4) Let's redefine the coordinate system, then... into axes D₁ and D₂. Notice that there is more spread among D₁ than D₂. PCA will recognize that and automatically make D₁ the first dimension, and D₂ the second dimension.

In summary, PCA can:

→ Transform data

- Removes correlation between data
- Ranks coordinates by importance

→ Concentrate on the first n principal components

- Reduces effect on randomness
- Earlier principal components are likely to have higher signal-to-noise ratio

... Basically,

you start with a large number of factors, then use PCA to generate a small number of factors that you can use to fit the model.

Math of PCA

To go through math of PCA, you first need to understand eigenvalues/eigen vectors.

A refresher:

ex) Linear Algebra:

$$\begin{pmatrix} 1 & 2 & 0 \\ 2 & 1 & 0 \\ 0 & 0 & -3 \end{pmatrix} \times \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 3 \\ 3 \\ 0 \end{pmatrix} = 3 \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}$$

eigenvalue (scalar value)

eigenvector (vector value)

⇒ (magnitude)(direction)

⇒ (eigenvalue)(eigenvector)

i.e. $\Rightarrow 3 \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}$

Avec transformation $A\vec{v} \Rightarrow 3\vec{v}$

eigenvalue

eigenvector

vectors don't change directions, just magnitude.

Say we have a square matrix.....

- \vec{v} is a vector such that: $A\vec{v} = \lambda\vec{v}$
- \vec{v} is an eigenvector of A
- λ is an eigenvalue of A: $\det(A - \lambda I) = 0$

Given λ , solve

$A\vec{v} = \lambda\vec{v}$ to find eigenvalue \vec{v} .

⇒ Software can do this

Now, back to PCA math.

Essentially, PCA uses the properties of eigenvectors and eigenvalues to get not just transformed set of coordinated directions, but also for the directions to be all orthogonal to each other.

Scaled matrix X of data:

- X : initial matrix of data
- X_{ij} = jth factor value for data point i, after scaling
- Scale such that $\frac{1}{m} \sum_i X_{ij} = \mu_j = 0$

Need to find eigenvectors v_1, \dots, v_n of $(X^T X)$ { $X^T X$ transposed $\times \times$ }

- V : matrix of eigenvectors (sorted by smallest to largest eigenvalue)
- $V = [v_1 \ v_2 \ \dots]$, where v_j is the jth eigenvector of $X^T X$

Find the principal components....

- Multiply X by the eigenvectors
 - ↳ First component is Xv_1 , second component is Xv_2 , etc.
- Xv_1, Xv_2, \dots, Xv_n are the principal components
 - ↳ these are the set of orthogonal coordinate directions
- If a lot of data and just trying to remove correlation: use all PCs.
But if you're trying to find a smaller number of variables in model, then pick n number of PCs for model.

This is all assuming linear transformation.

If nonlinear functions \rightarrow use kernels.

Week 7: Advanced Regression

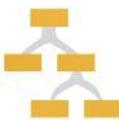
Sunday, March 6, 2022 3:15 PM

We've built, fitted, and evaluated regression models.

Now, let's divide data sets and have different models for each data subset.

We can do this with trees.

Classification and Regression Trees (CART)



This is a decision tree.

They can help us make a decision, and we can separate our models into different "systems" to our liking.

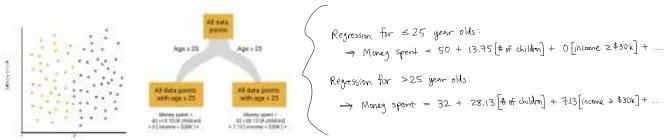
Example for a single regression model

Say we want to find the impact of a marketing email on response spending. Predictors could include:

- Demographic: age, sex, # children, income, etc.
- Marketing factors: avg spent per month in site
- Binary factors (was the email received? yes/no)

But what if some of those factors behave differently in different combinations?
i.e., maybe the coefficient for receiving a marketing email should really differ by age.

We can split our model into two branches....



It doesn't end here!

We can make as many branches as we want. When a branch "ends," they're called leaves.

For each leaf, we'll run a separate regression to find the individual set of coefficients. { only do regression on the LEAF/leaves }

The tree approach can:

- Be used as a flow chart to determine which model to use
- Can compare one leaf to another to see if the model needs to be improved, i.e. can compare R² values for each leaf

Note we can use decision trees for any model, not just regression!

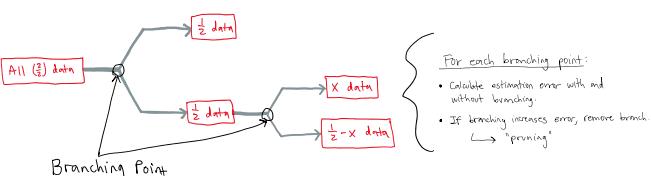
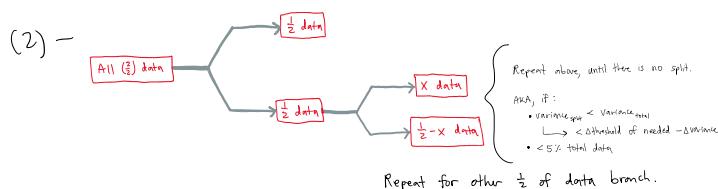
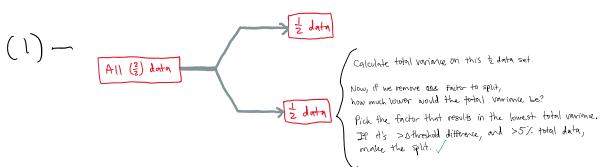
So, how do we branch, and when do we stop branching?

In theory, we could have a branch for every one data point in it....
But that would be ridiculous and a very poor predictor.

Rule of thumb: each leaf contains at least 5% of the data

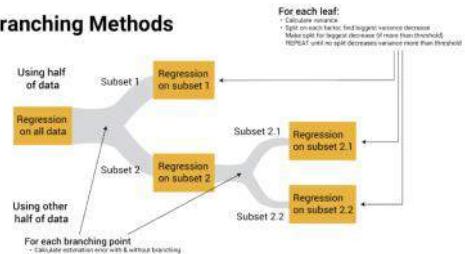
If the branch can lead to overfitting, reject it.

To begin branching, start with one factor at a time.



Lecture's Flow Chart

Branching Methods



Random Forests

.... Forests = many trees together

Random forests introduce randomness and generate many different trees with varying strengths and weaknesses.
Idea: overall, the average of all these trees is better than a single tree with specific strengths and weaknesses. → helps prevent overfitting

Steps to introduce randomness:

- (1) give each tree slightly different datasets ('bootstrapping' process)
 - ↳ tree could have multiple of the same points, and/or none of others

- (2) randomly branch
 - ↳ randomly choose a small number of features X (instead of doing 1 feature at a time), and then choose the best feature with X to branch on
 - a common number of features to use is $1 + \log(n)$

- (3) don't prune the tree
 - ↳ keep all the branches, even if error is high

We may end up with 500-1000 trees.... which one do we use?
 ↳ It depends.

- Regression Trees ⇒ use the average predicted response
- Classification Tree ⇒ use the mode, the most common predicted response

Unfortunately, random forests don't give us a specific regression/classification model from the data.

Still, used as a "default" model → if you want quick results and don't have a reason to try another method, use random forests.

But you may want to use another model if you need detailed insights to what's going on.

Confusion Matrices

Can measure how well a classification-type model works, since it can incorrectly/correctly classify something.

		Model's Classification	
		Yes	No
True Classification	Yes	Correct	Incorrect
	No	Incorrect	Correct

		Model's Classification	
		Yes	No
True Classification	Yes	True Positive (TP)	False Negative (FN)
	No	False Positive (FP)	True Negative (TN)

Guidelines

- True — model got it right
- False — model got it wrong
- Positive — model says it's in the category
- Negative — model says it's NOT in the category

Let's see an example.

ex) Say there is an email spam filter, and 50% of the incoming email is spam

Confusion matrix:

		Model's Classification	
		Real	Spam
True Classification	Real	True Positive (TP) 490	False Negative (FN) 10
	Spam	False Positive (FP) 100	True Negative (TN) 400

Question:

- a) What fraction of the inbox will be spam with this filter?

$$\Rightarrow \text{total emails in inbox} = TP + FP \Rightarrow 490 + 100$$

$$\hookrightarrow \text{fraction of spam in inbox} = \frac{FP}{TP+FP} = \frac{100}{490+100} \approx 17\% \rightarrow 1 \text{ out of } 6$$

- b) What fraction of the inbox is real email flagged as spam?

$$\Rightarrow \text{total emails in inbox} = TP + FP \Rightarrow 490 + 100$$

$$\hookrightarrow \text{fraction of lost emails} = \frac{FN}{TP+FP} = \frac{10}{490+100} \approx 2\%$$

Some common equations ...

$$\text{Sensitivity} = \frac{TP}{TP + FN} \Rightarrow \text{fraction of category members correctly identified}$$

}

$$\text{Specificity} = \frac{TN}{TN + FP} \Rightarrow \text{fraction of non-category members correctly identified}$$

Others that you don't need to memorize:



Evaluating a model's quality using a confusion matrix

Let's go back to the spam example.

Question :

Evaluating a model's quality

Let's go back to the spam example.

Model's Classification			
Real	Spam		
True Classification	Real	True Positive (TP) 490	False Negative (FN) 10
	Spam	False Positive (FP) 100	True Negative (TN) 400

Question:

what is the cost of lost productivity [reading spam]?

Given:

- \$0 lost for correct classifications $\rightarrow TN, TP$
- \$0.04 lost to read spam $\rightarrow FP$
- \$1 lost when missing a real message $\rightarrow FN$

a) If 50% of email is spam:

$$\begin{aligned} \Rightarrow \text{Total cost lost} &= TP(\$0) + FN(\$1) + FP(\$0.04) + TN(\$0) \\ &= (490)(\$0) + (10)(\$1) + (100)(\$0.04) + (400)(\$0) \\ &= \$14 \\ &= \$0.014/\text{email} = 1.4\text{¢}/\text{email} \end{aligned}$$

a) If 40% of email is spam:

\hookleftarrow problem: given data is for 50% rate.....

so, we need to scale our math.

.... 40% is spam, 60% is real.....

Spam gets scaled down from 50% to 40%;
Real email gets scaled up from 50% to 60%.

$$\begin{aligned} \Rightarrow \text{Total cost lost} &= (490)\left(\frac{0.4}{0.5}\right)(\$0) + (10)\left(\frac{0.6}{0.5}\right)(\$1) + (100)\left(\frac{0.4}{0.5}\right)(\$0.04) + (400)\left(\frac{0.6}{0.5}\right)(\$0) \\ &= \$15.2 \\ &= 15.2\text{¢}/\text{email} \end{aligned}$$

Why is this useful?

So we reduce the original 17% inbox spam to 10%,
but in return, it means there are more FNs (real emails flagged as spam).

Is this really better?

With the equations above, we can calculate a higher cost of productivity lost.

\curvearrowleft This is a great example of weighing different attributes. (see SVM notes)

Logistics Regression

We've worked with models that are continuous....

But what if we want to estimate a probability? \Rightarrow use logistic regression model

- ex) • probability that the loan recipient will repay the entire loan in time \rightarrow answer is 0 through 1 \Rightarrow logistic regression
 • yes or no: the recipient will pay back the loan in time \rightarrow answer is 0 or 1 \Rightarrow linear regression

Basically,
Linear regression is a straight line.
Logistic regression is a curved line.

\Rightarrow Standard Linear Regression Model:

$$y = a_0 + a_1 x_1 + a_2 x_2 + \dots + a_j x_j$$

\Rightarrow Logistic Regression Model:

$$P = \frac{1}{1 + \exp[-(a_0 + a_1 x_1 + a_2 x_2 + \dots + a_j x_j)]}$$

$$\log\left(\frac{P}{1-P}\right) = a_0 + a_1 x_1 + a_2 x_2 + \dots + a_j x_j$$

where:

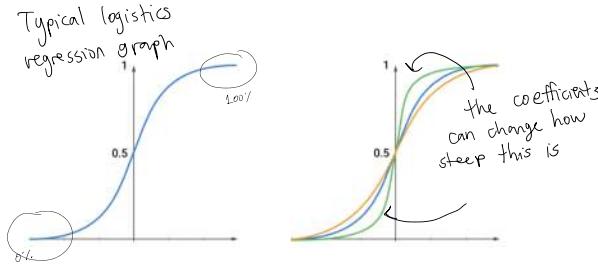
P = probability of event you want to observe

If:

$$\begin{cases} a_0 + a_1 x_1 + a_2 x_2 + \dots + a_j x_j = -\infty \Rightarrow P = 0 \\ a_0 + a_1 x_1 + a_2 x_2 + \dots + a_j x_j = +\infty \Rightarrow P = 1 \end{cases}$$

Typical logistics
regression graph





Logistic Regression vs. Linear Regression

Similarities:

- Transformations of input data
- Consider interaction terms
- Variable Selection
- Logistic regression trees
- Random logistic regression forests

Differences:

- (Relatively) Longer to calculate
- No closed-form solution
- Understanding model quality

Linear regression has an R^2 value (= fraction of variance explained by model)

.... Logistics regression doesn't have one.

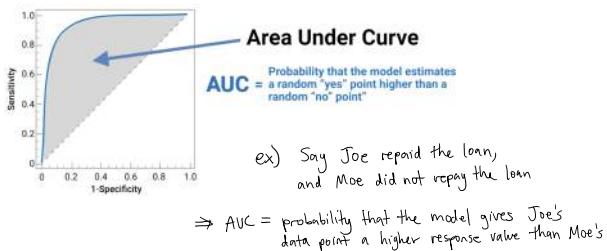
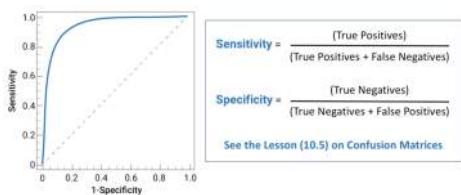
But, people made one up that sort of works, called the "pseudo R^2 value" to compare different models
Just keep in mind that it doesn't measure the same thing (variance)

Thresholding

Answer "yes" if probability p is at least some number; otherwise, answer "no."
i.e. If $p \geq 0.7$, give loan. Otherwise, no.

A visual way of looking at this is the Receiver Operating Characteristic (ROC) Curve.

It plots the (1-specificity) and (sensitivity):



For reference: $AUC = 0.5 \Rightarrow$ basically just guessing

Advanced Topics in Regression

- Poisson Regression
Use when response follows a Poisson distribution: $f(z) = \frac{\lambda^z e^{-\lambda}}{z!}$
- Regression Spline
Use when you want to fit different functions to different parts of the data set; smooth connections between points
- Bayesian Regression
Use when combining expert opinion/normal distribution with small amount of data

- **Regression Spline**
Use when you want to fit different functions to different parts of the data set; smooth connections between points
- **Bayesian Regression**
Use when combining expert opinion/prior distribution with small amount of data
- **k-Nearest Neighbor (KNN) Regression**
Use when needing to predict a response; use average response of k closest data points

$$(SSE) \text{ Minimize } \sum_{i=1}^n [y_i - (a_0 + a_1x_{1i} + a_2x_{2i} + \dots + a_jx_{ji})]^2$$

$$\Rightarrow \text{Subject to } \sum_{i=1}^j |a_i| \leq \tau \rightarrow \text{restriction: sum of coefficients can't be too large}$$

giving the regression a budget to use on coefficients... use budget to prioritize important coefficients... unimportant coefficients will be less, so they won't be used in the model.

How to pick τ ? Depends on:

- (1) Number of variables
- (2) Quality of model as you allow more variables

Use Lasso approach with many varying τ values to see which gives you the best trade-off

(5) Elastic Net: same as LASSO, but constrain a combination of the absolute values of the coefficients and their squares

$$\text{Minimize } \sum_{i=1}^n [y_i - (a_0 + a_1x_{1i} + a_2x_{2i} + \dots + a_jx_{ji})]^2$$

$$\Rightarrow \text{Subject to } \lambda \sum_{i=1}^j |a_i| + (1-\lambda) \sum_{i=1}^j a_i^2 \leq \tau$$

NOT FOR VARIABLE SELECTION

(6) Ridge Regression: same as Elastic Net, but no absolute value term / $\lambda = 0$

$$\text{Minimize } \sum_{i=1}^n [y_i - (a_0 + a_1x_{1i} + a_2x_{2i} + \dots + a_jx_{ji})]^2$$

$$\Rightarrow \text{Subject to } \sum_{i=1}^j a_i^2 \leq \tau$$

So, now we have 3 greedy approaches to variable selection, and 3 global ones.

Which method do we pick?

Pro/con comparison:

- | | | |
|-----------------------|--|--|
| <u>Greedy Methods</u> | $\left\{ \begin{array}{l} \cdot \text{Forward Selection} \\ \cdot \text{Backwards Elimination} \\ \cdot \text{Stepwise Regression} \end{array} \right\}$ | $\left\{ \begin{array}{l} \cdot \text{Quick & fast} \\ \cdot \text{Good for initial data analysis} \\ \cdot \text{Dosen't check for random effects as well; not as thorough} \\ \cdot \text{Thus, won't perform as well when tested on other datasets} \end{array} \right\}$ |
|-----------------------|--|--|
- ↳ more common among the three

- | | | |
|-----------------------|--|--|
| <u>Global Methods</u> | $\left\{ \begin{array}{l} \cdot \text{LASSO} \\ \cdot \text{Elastic net} \end{array} \right\}$ | $\left\{ \begin{array}{l} \cdot \text{More advanced} \\ \cdot \text{Better predictive models} \\ \cdot \text{Slower} \end{array} \right\}$ |
|-----------------------|--|--|

Basically:

- If introducing data exploration: use greedy (step-wise) methods first, then build a more refined model later using LASSO/Elastic Net.
- If not introducing data: use LASSO or Elastic Net

What's the difference between LASSO & Elastic Net? ...

Regularized Regression

Lasso: (Some coefficients forced to 0 to simplify model)

$$\bullet \text{ Minimize } \sum_{i=1}^n (y_i - (a_0 + a_1x_{1i} + a_2x_{2i} + \dots + a_jx_{ji}))^2$$

Subject to $\sum_{i=1}^j |a_i| \leq \tau$

(LASSO): Helps denote the number of factors with non-zero coefficients

picks important predictors and discards the rest

→ has zero-value coefficients

Ridge: (Coefficients shrink toward 0 to reduce variance in estimate)

$$\bullet \text{ Minimize } \sum_{i=1}^n (y_i - (a_0 + a_1x_{1i} + a_2x_{2i} + \dots + a_jx_{ji}))^2$$

Subject to $\sum_{i=1}^j a_i^2 \leq \tau$

(Ridge): Quadratic

Pushes coefficients towards zero, but not equal to zero

shrink the coefficients of correlated predictors towards each other

↳ low variance ... but trading off loss for variance = better predictive models

(Cherry variance vs. bias) ... so ridge model = family / drama!

Elastic net

$$\bullet \text{ Minimize } \sum_{i=1}^n (y_i - (a_0 + a_1x_{1i} + a_2x_{2i} + \dots + a_jx_{ji}))^2$$

Subject to $\lambda \sum_{i=1}^j |a_i| + (1-\lambda) \sum_{i=1}^j a_i^2 \leq \tau$

Elastic Net = LASSO + Ridge

↳ variable selection benefits
• exactly pushes some coefficients to zero, even if irrelevant.
i.e. if hard to choose between 2 very similar values, it'll push one and discard the other

i.e. 2 additive terms but one is expensive to compute, so we can choose the expensive one unless you hard set and usually choose can better fit the model

Analytics is an art.

Makes provide insight and direction, but human intuition is also beneficial.

→ Try all the methods, and use your best judgement when comparing results.

So far for this course, the data is given or is easily collected.

But that's not always the case.

Sometimes we don't have data, or a full data set is hard to get.

i.e. Data can come from a survey, but surveys complete in the word is expensive and impossible.

So, instead, like last page, we can represent things.

→ This is easier and less expensive.

Using cars (size of 2 cm³), and car is very low price.

Design of Experiments (DOE):

designing a way to collect the basic kind of data that answers our questions, quickly and efficiently.

Concepts:

 (1) Comparison: compare car [below] to the car [above] → compares one factor at a time.

(2) Control: making sure all other factors outside the comparison are the same.

(3) Blocking: a blocking factor is something that could create variation.

→ Type of car is a blocking factor because it's part of the difference between the two cars, so we want to make sure that's not a confounding factor.

A few different types of DOE....

A/B Testing: a test to see which of the two alternatives is better

When choosing between 2 alternatives, can use A/B testing as long as those are the true.

1. Collect data quickly enough to get an answer in time.
2. Data must come from a representative sample of the total population.

3. Amount of data is small compared to the total population.

→ 20 visitors per day over the last 200 days would give us 4000 visitors.

4. Which of the two ads are more effective?



Suppose we've been ads for the first 2000 days. The ad is shown half the time, and it's shown to 10% of the people who look at the ad.

Now we're going to show ad A to 10% of the people who look at the ad.

Sample 1

Sample 2

Hypothesis test:

B is significantly better.

Hypothesis test:

No significant difference.

But what if we have several alternatives, not just two?

- If different factors/differences involved, we want to know why A is better than B?

Factorial Testing: a test to see which of many alternatives is better

Break the comparisons into factors rather than an absolute A and B.

... If you're in the position to decide what data to collect and what comparisons to try, factorial design can be a powerful tool.

5. Suppose ad B generates 50% more clicks, but why?



Different font, different reading, different background color...

There are all factors we can test.

Full factorial Design: test the effectiveness of every combination

Fractional Factorial Design: test the effectiveness of a subset of combinations

→ estimates all effects by comparing some conditions

		(A)	
		Background	
Font	Reading	Background	
Arial	MIS in Analytics	White	Gold
Arial	MIS in Analytics	Gold	White
Arial	Master of Science in Analytics	White	Gold
Arial	Master of Science in Analytics	Gold	White
Roboto	MIS in Analytics	White	Gold
Roboto	MIS in Analytics	Gold	White
Roboto	Master of Science in Analytics	White	Gold
Roboto	Master of Science in Analytics	Gold	White

→ Full factorial design: test every combination of the ad

- 2 fonts × 2 reading × 2 backgrounds = 8 combinations
- we want to determine influence of each factor

→ But what if there are many combinations?

- 7 factors, 3 levels each = $3^7 = 2187$ combinations
- ... That's way too many! → use fractional factorial design instead

→ Fractional factorial design: test subset of combinations

→ Test each since the same # of times

Balance Design: → Test each pair the same # of times

Still, there is a trade-off:

- More Information vs. Immediate Value

→ Explanation vs. Explanation

This trade-off problem is called the Multi-Armed Bandit.

- ... name comes from gambling slot machines. To play the highest payout machine, the player would need to hit all of them.

Iteration Process:

- Start with no information
 - Equal probability of selecting each alternative
- Perform some tests
 - Gather information
- Based on new information, choose next one on which one you think/know is the best
 - Update Information

Parameters in the Iteration:

- Number of trials between consecutive probabilities

- How to update the probabilities

Collect data.

Update possibilities.

Continue until you've found the best one.

Very good for learning faster on the fly and creating more value along the way.

(Noteworthy!) Better than running a fixed, large number of tests.

Probability Distributions

→ Software can determine if output fits probability distribution

Sometimes you don't need the complexity of factors in a model.

Sometimes a simple model is all you need!

→ Probability distributions can form the backbone of such simple models → it's the simple approach.... i.e. Bernoulli example is

Eespecially true when the only data we have is the response, and/or

it would be hard to collect/analyze additional information.

→ Bernoulli Distribution: probability of yes/no (success/fail) with the outcome of the single trial of the event → i.e. coin flip

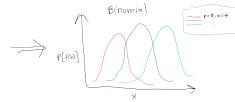
→ consists of one trial

$$\text{Probability Mass Function: } \begin{cases} P(X=1) = p \\ P(X=0) = 1-p \end{cases}$$

Independent → **Binomial Distribution:** probability of yes/no (success/fail) with the outcome of the multiple trials of the event → i.e. count of successful coin flips in a trial

→ consists of many Bernoulli(p) trials → probability of x successes out of n Bernoulli(p) trials

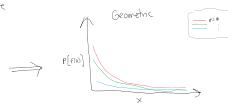
$$\text{Probability Mass Function: } \left\{ \begin{array}{l} P(X=x) = \binom{n}{x} p^x (1-p)^{n-x} \\ = \frac{n!}{(n-x)! x!} p^x (1-p)^{n-x} \end{array} \right. \Rightarrow \text{As } n, \text{ binomial distributions} \rightarrow \text{normal distribution}$$



→ **Geometric Distribution:** probability of # of failures before the first success → i.e. unsuccessful coin flips before first success OR fail-to-success ratio of door-to-door sales

→ consists of many Bernoulli(p) trials → probability of having x Bernoulli(p) failures until the first success
OR probability of having x Bernoulli($1-p$) successes until the first failure

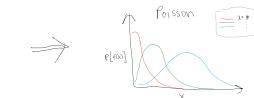
$$\text{Probability Mass Function: } \left\{ \begin{array}{l} \text{failures until first success} \\ P(X=x) = (1-p)^x p \\ \text{successes until first failure} \\ P(X=x) = (1-p)^x \end{array} \right. \Rightarrow \text{Can also use geometric distribution to check if factors are independent and identically distributed (iid)}$$



Independent

→ **Poisson Distribution:** probability of having a certain # of events occurring in a fixed time interval → i.e. # of customers arriving in a restaurant

$$\text{Probability Mass Function: } \left\{ f_x(x) = \frac{\lambda^x e^{-\lambda}}{x!} \right. \left. \begin{array}{l} \cdot \text{Good at modeling random arrivals} \\ \cdot \lambda: \text{average number of arrivals per time period} \\ \cdot \text{Arrivals are independent and identically distributed (iid)} \end{array} \right.$$



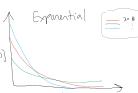
If Poisson, then it's "Memoryless": distribution of time to next interval = initial distribution of time to next arrival
↳ same as last pointer, what happened in the past, all that matters is where we are now

↳ if not memoryless, then NOT Poisson

Independent

→ **Exponential Distribution:** probability of time waiting before an event occurs → i.e. time between customers arriving in a restaurant on a given day

$$\text{Probability Mass Function: } \left\{ f_x(x) = \lambda e^{-\lambda x} \right. \left. \begin{array}{l} \text{If arrivals are Poisson}(\lambda) \\ \text{then the time between successive arrivals is exponential}(\lambda) \text{ distribution.} \end{array} \right.$$



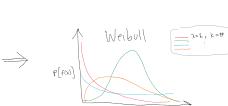
If exponential then it's "Memoryless": distribution of remaining time = initial distribution of time
(λ stays the same)

↳ If not memoryless, then NOT exponential

→ **Weibull Distribution:** probability of the time in which something fails → i.e. lifetime curve of plane battery before recommended to be replaced

→ Weibull vs. Geometric
(time between failures) vs. (# of times until failures)

$$\text{Probability Mass Function: } \left\{ f_x(x) = \frac{k}{\lambda} \left(\frac{x}{\lambda} \right)^{k-1} e^{-\left(\frac{x}{\lambda} \right)^k} \right. \left. \begin{array}{l} \lambda: \text{Scale Parameter} \\ k: \text{shape parameter} \end{array} \right.$$



• $k > 1$ → modeling when failure rate decreases with time, "most things fail first" i.e. parts with defects

• $k < 1$ → modeling when failure rate increases with time, "things that never end" i.e. tires

• $k = 1$ → modeling when failure rate is constant with time...
if $k = 1$, then it becomes an exponential distribution.

$$f_x(x) = \left(\frac{1}{\lambda} \right) \left(\frac{x}{\lambda} \right)^{k-1} e^{-\left(\frac{x}{\lambda} \right)^k}$$

$$= \frac{1}{\lambda} e^{-\frac{x}{\lambda}} = \lambda e^{-\lambda x} \quad (\text{exponential case})$$

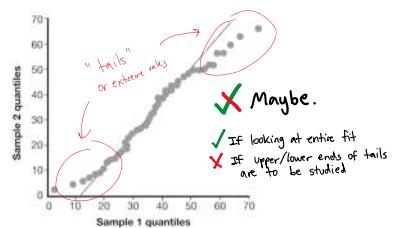
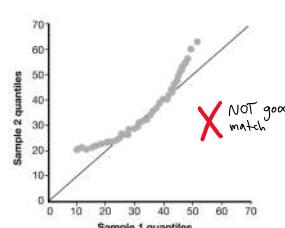
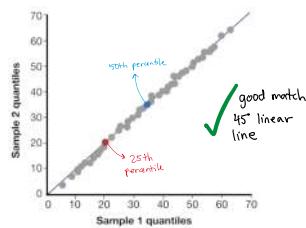
Q-Q Plots: visualization on whether two distributions of data are about the same, or whether one data set is distributed similarly to a probability distribution.

Q-Q plots plot percentile/quantile of data sets to show that they are/are not distributed the same

The idea is that, despite any variations, two similar distributions should have about the same value at each quantile

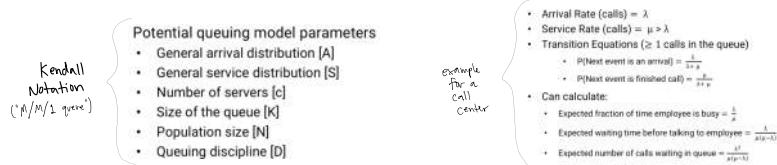
Plotting (x, y), where:

$x = \text{Sample 1's } z\text{-th percentile value}$
 $y = \text{Sample 2's } z\text{-th percentile value}$



Alternatively, x-axis is data
& y-axis is ideal/theoretical values for perfect distribution.

Queuing Models: study of waiting lines / queues.



Simulation: build a model to watch its behavior

Pick 1

Deterministic Simulations: no randomness; never vary; same inputs results in same outputs

Stochastic Simulations: has randomness; varies; different inputs results in different outputs

Continuous-Time Simulations: changes happen continually (often modeled with differential equations)

i.e. chemical processes, propagation of diseases

Pick 1

Discrete Event Simulations: changes happen at discrete time points

i.e. call center simulations, airport departure, car wash

Focus on lesson: discrete-event stochastic simulations

- Valuable when systems have high variability/randomness
- Using average values isn't good enough
 - Use a lot of replications (number of runs)!
 - 1 replication = 1 data point (may be unrepresentative)

Simulation Validation

- Use real data to validate your simulation is giving reasonable results
 - Real and simulated averages don't match \rightarrow problem
 - Averages match, variances don't match \rightarrow problem

We have learned about descriptive & predictive analytics.

When talking about simulation, we can now discuss **prescriptive analytics**, aka the "what if?" & "what should I do?" questions

- Sample:
- How are the capital resources throughout the life of the company invested?
 - How valuable would hiring an extra call center worker be?
 - What would be the impact of adding a new gate at an airport - one for each gate, one for every five gates, or a floating pool at each half side of the terminal?

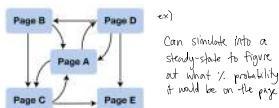
Simulation can be a powerful prescriptive tool.
But remember, it's only as good as its input & assumptions.

Simulation

- Simulation software
 - Elements of model include
 - Entities that move through simulation (e.g., bags, people, etc.)
 - Modules parts of process (e.g., queues, storage, etc.)
 - Actors
 - Decisions (e.g., waitlist)
 - Decide points
 - Statistical tracking
 - Often "drag and drop" programming style
 - Under the hood
 - (Process) function, marker generation, tracking variables, etc.

Markov Chains: stochastic model describing a sequence of possible memoryless events

- Memoryless: state transitions only depend on the most recent state
- Most systems don't exhibit this property, but if they do, it's a powerful tool



Week 10: Missing Data, Optimization

Saturday, April 9, 2022 10:28 PM

So far, we assume we have data. Sometimes it's big data—billions of data points.

But sometimes, not all of those data points are available.

Sometimes there is missing data.

Car #	Color	Style	Miles	Price
1	White	Sedan	20,000	\$12K
2	White	Coupe	30,000	\$12K
3	Black	Sedan	17,000	\$10K
4	Gray	Coupe	24,000	\$11K
5	Red	SUV	29,000	\$20K
6	White	Pickup	24,000	\$13K
7	Gray	Sedan	21,000	\$13K

Can just be missing or wrong, but also there could be a pattern.

i.e. Missing data may occur during:

- Temperature Readings → gauges break down
- Filling out forms → manually filled out, possibly left blank, unanswerable data
- Income Survey → higher income people less likely to report income ↙ **BIAS**
- Radar Gun → may treat slow cars as moving and not show a speed

What to do with missing data?

- (1) Throw away
- (2) Use categorical variables ("N/A", "missing") to indicate missing data
- (3) Impute/estimate data

(1) Throw Away

Pros { Not potentially introducing errors
Easy to implement

Cons { Don't want to lose too many data points
Potential for censored or biased missing data → Missing not at random = biased data

(2) Categorical Variable Approach

→ Categorical variable has missing data:

Add category: missing

Car #	Color	Style
1	White	Sedan
2	White	Coupe
3	Black	Sedan
4	Gray	Coupe
5	Red	SUV
6	White	Pickup
7	Gray	Sedan

Car #	Color	Style
1	White	Sedan
2	White	Coupe
3	X	Sedan
4	X	Coupe
5	X	SUV
6	X	Pickup
7	X	Sedan

→ Quantitative variable has missing data:

- All missing values = 0
- Add new categorical variable

Car #	Var 1	Var 2	Var 3	Var 4
1	1	2	3	4
2	1	2	3	4
3	1	2	3	4
4	1	2	3	4
5	1	2	3	4
6	1	2	3	4
7	1	2	3	4

Car #	Var 1	Var 2	Var 3	Var 4	Var 5
1	1	2	3	4	0
2	1	2	3	4	0
3	X	X	X	X	0
4	X	X	X	X	0
5	X	X	X	X	0
6	X	X	X	X	0
7	X	X	X	X	0

Car #	Var 1	Var 2	Var 3	Var 4	Var 5	Var 6
1	1	2	3	4	0	0
2	1	2	3	4	0	0
3	X	X	X	X	0	0
4	X	X	X	X	0	0
5	X	X	X	X	0	0
6	X	X	X	X	0	0
7	X	X	X	X	0	0

(3) Imputation

No more than 5% of data should be imputed to avoid overfitting (data used twice)

3 methods:

1. Midrange Value

Mean, Median (numeric), Mode (categorical)

Advantages:

- Hedge against being too wrong
- Easy to compute

Disadvantages:

- Imputation can be biased
i.e. mean/mode will be underestimated in income survey
since people with high income are less likely to answer

2. Use a Predictive Model

Predicting the value of missing data based on the other variables.
...like regression.

Advantages:

- Reduces/Eliminates bias

Disadvantages:

- Complex: build fit, validate, & test to estimate missing value
- Overfitting: using the data twice; once for imputation, second to fit model
- Imputation does not capture all variability

3. Perturbation

+/- random amount from a normally-distributed distribution
... 'noise' or random value

Advantages:

- More accurate variability (Artificially bring in variability)

Disadvantages:

- Less accurate on average. (better off using predictive model)



Question: Does imputation introduce additional errors?

- Imputation Error + Perturbation Error (if applicable) + Model Error
 - It's very hard to determine each error component, and how they add up.
 - We can use a test data set to evaluate the overall model's performance.
 - That's the most ideal solution, despite the fact that the test data is missing some data
- So, yes, errors do get added.
- That doesn't mean imputation is inferior to using regular data.
 - Regular data can be just as worrisome
 - It's almost impossible to have a dataset where each value is actually correct.

Data is always imperfect:

- Errors
- Outliers
- Missing data

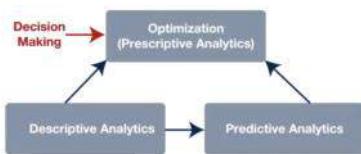
Optimization

clustering, regression, variable selection, time series...

Optimization is used by many statistical models to get answers in descriptive and predictive analytics.

But optimization can also be used for prescription analytics: "given what I know and what I predict, what's the best course of action to take?"

- Optimization can answer questions like:
- Which airplane flights should be scheduled for each flight over the next week to meet the expected maintenance requirements at Other Airports at lowest cost, making sure we don't violate any federal regulations or union contract requirements and addressing for unexpected breakdowns?
 - How much crude oil should be sent by tank and by pipeline from each of four refineries to each of three distribution centers while keeping costs low?
 - What web pages should be optimized in a search filter, and how many copies of each should be stored to maximize the profit made by responding quickly to visitors?
 - How should a large machine shop reorganize its production to get maximum throughput while meeting all clients' deadlines, paying into account the many that come up on fuel inspection and need to be remade?
 - Where the source route flies my plane to the airport, given current and predicted traffic?



Predictive/Descriptive Models:
Software automates solution
Software automates model building

Optimization/Prescriptive Models:
Software automates solution
YOU do the model building → building models requires human insight & expertise

Optimization models have three main components:

1) Variables

Decisions to be made.

Example political candidate scheduling
 x_i = total time spent in state i
 y_i = number of visits to state i
 $z_i = 1$ if state i is ever visited, 0 if not
 $w_{id} =$ time spent in state i on day d
 $v_{id} = 1$ if state i visited on day d , 0 if not

2) Constraints

Restrictions on variable values.

Constraints:
 $\sum x_i \leq 30$
At least 3 Florida visits in days 24-30
 $\sum_{d=24}^{30} v_{i,d} \geq 3$
Total visits must add up correctly
 $\sum_d v_{id} = y_i$

3) Objective Function

Solution quality measure; what are you trying to minimize/maximize?

Can rank/assess variables given restrictions.

$$\sum_i \left(\alpha p_i \sqrt{x_i} + \frac{1}{3} \sum_j x_j + \beta v_{id} f_d \right)$$

Determined by other analytics models!

Pretty common for optimization models to require input that's found as the output from other models

Some more definitions...

Solution: values for each variable

Feasible Solution: variable values that satisfy all constraints

Optimal Solution: feasible solution with the best objective value

Modeling is an Art: Examples

Diet Problem (US Army):

Diet Problem (US Army):

Satisfy soldiers' nutritional requirements at minimum cost

Given :

n foods
 m nutrients
 a_{ij} = amount of nutrient j per unit of food i
 m_j = minimum daily intake of nutrient j
 M_j = maximum daily intake of nutrient j
 c_i = per-unit cost of food i

Variables : x_i = amount of food i in daily diet

Constraints : $\sum_j a_{ij}x_i \geq m_j$ for each nutrient j
 $\sum_j a_{ij}x_i \leq M_j$ for each nutrient j
 $x_i \geq 0$ for each food i

Objective Function : Minimize $\sum_i c_i x_i$

✓ Optimization Model.

Other complexities to consider:

- Variety in diets
- Seasonal cost variation
- Taste of food
- Taste of combinations of foods

Example 2

- Call center scheduling
- Meet forecasted demand d_i for each day of the week i
- Workers work 5 days in a row, then 2 days off
- Minimize worker-days used

~~Variables~~
 x_i = number of people working on day i

~~Objective function~~
minimize $x_{\text{Sunday}} + x_{\text{Monday}} + \dots + x_{\text{Saturday}}$

~~Constraints~~
Meet demand: $x_i = d_i$ for all days i
Integer workers: x_i is integer for all days i
5-day requirement: ???

Variables
 x_i = number of people who start working on day i

Objective function
minimize $5(x_{\text{Sunday}} + x_{\text{Monday}} + \dots + x_{\text{Saturday}})$

Constraints
Meet demand: \sum_j working on day i $x_j \geq d_i$
Ex: $x_{\text{Fr}} + x_{\text{Sat}} + x_{\text{Sun}} + x_{\text{Mon}} + x_{\text{Tue}} \geq d_{\text{Tue}}$

Non-negativity: $x_i \geq 0$ for all days i
Integerality: x_i is integer for all days i

Modeling with Binary Variables :

Example: Stock Market Investment

Invest to balance return and risk

B = investment budget

n = number of stocks available

r_i = expected return of stock i relative to market

Q_{ij} = covariance of returns of stocks i and j

m_i = minimum dollar amount for each stock

Personal constraints (more examples)

If invest in energy, invest in at least 5

Option 2: $z_{\text{Energy}} = 1$ if invest in energy, 0 if not

$$\sum_{j \in \text{energy}} y_j \geq 5z_{\text{Energy}}$$

$$z_{\text{Energy}} \geq y_j \text{ for all energy stocks } i$$

Variables

x_i = amount invested in stock i
 $y_i = 1$ if invest in stock i , 0 if not

Constraints

$$\sum_i x_i \leq B$$

$$x_i \geq 0 \text{ for all stocks } i$$

$$x_i \leq B y_i \text{ for all stocks } i$$

$$x_i \geq m_i y_i \text{ for all stocks } i$$

$$y_{\text{Tesla}} = 1$$

$$y_{\text{Amazon}} + y_{\text{Google}} + y_{\text{Apple}} \geq 1$$

$$y_{\text{FedEx}} = y_{\text{UPS}}$$

$$y_{\text{Coca-Cola}} = 1 - y_{\text{PepsiCo}}$$

Objective function

$$\text{Maximize } \sum_i r_i x_i - \theta \sum_i \sum_j Q_{ij} x_i x_j$$

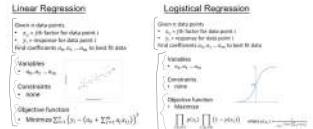
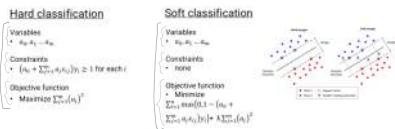
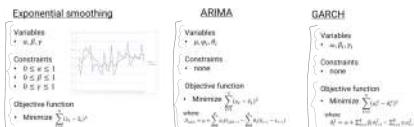
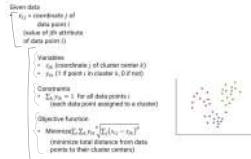
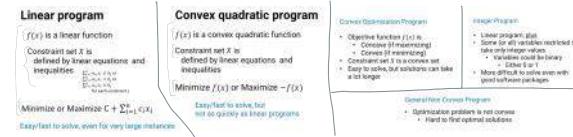
$$- \sum_i y_i$$

Optimization is important in analytics/statistics.

In fact, it's used in all of our modeling techniques thus far:

Here are examples.

First, a note...

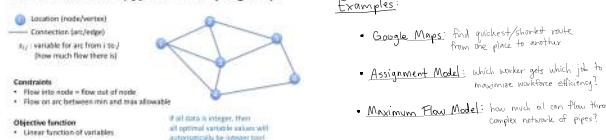
Notation confusion: statistics vs. optimization**Regression****Support vector machine models****Time series models****k-means for clustering****Classification of Optimization Models**

Objective Function	Constraints	Program Type	Difficulty to Solve	Time to Solve
Linear	Linear	Linear	Easy	Fast
Quadratic	Linear	Quadratic	Easy	Faster, but slower than Linear Objective Function
Quadratic	Quadratic	Convex	Easy	Can Take Long
Linear/Quadratic	Linear, but integer variables are linear (subset of linear)	Integer	Hard	Slow
Linear/Quadratic	Linear, but integer variables are integers (subset of integers)	Integer	Harder	Slower
	General Non-Convex	Convex	Hardest	Slowest

From quickest/easiest to slowest/hardest:

1. Linear Programs
2. Convex Quadratic Programs
3. Convex Programs
4. Integer Programs
5. General Non-Convex Programs

If program is too hard, use a heuristic.

Network Models (type of linear program)**Examples:**

- Google Maps: find quickest/shorest route from one place to another
- Assignment Model: which worker gets which job to maximize workforce efficiency!
- Maximum Flow Model: how much oil can flow through complex network of pipes?

Stochastic Modeling

A model that accounts for randomness or uncertainty.

Use f :

- data/parameter isn't exactly known
- forecast value isn't exactly known

Two Approaches:**1. Model Conservatively****Not Conservative****Conservative**

Call center worker constraint:

 $x_1 + x_2 + x_3 + x_4 + x_5 + x_6 \leq 10$ \rightarrow not a conservative constraint

1. Model Conservatively

Not Conservative:
 Call center worker constraint
 $x_{\text{off}} + x_{\text{idle}} + x_{\text{out}} + x_{\text{on}} \geq d_{\text{target}}$
 $d = \text{number of customers starting today off on day } i$
 $\& \text{ if no customers start}$

Conservative:
 Call center worker constraint
 $x_{\text{off}} + x_{\text{idle}} + x_{\text{out}} + x_{\text{on}} \geq d_{\text{target}} + \theta$ → OR, with a *slack constraint* (possibly loose constraint)

2. Scenario Modeling

Scenarios:
 Scenario 1: very small returning flight
 Scenario 2: one major bug 3 days after launch
 Scenario 3: two major, immediate bugs
 Scenario 4: no major bugs, but demand after 10,000 airways

Robust Solution:
 model each scenario that satisfies every scenario constraint

Programming Models

The models we have seen so far are **mathematical programming models**.

As in, they have variables, constraints, and an objective function.

But other models may have a different structure:

Dynamic Program

- Assumes there is **no uncertainty**
- Relies on **states** (the exact situation and their values) and **decisions** (choices of next state).
- Systems are divided into states. The state specifies exactly what's going on in the system.
- Bellman's Equation: determine optimal decisions

Stochastic Dynamic Program

- Introduces uncertainty
- Basically dynamic programming model, but decisions have probabilities of next state

Markov Decision Process

- Stochastic dynamic program with discrete states and decisions
- Probabilities depend only on current state/decision.

Solving Optimization Models

We've learned how to build optimization models. Now we'll learn how to solve them.

Most optimization algorithms have these main steps:

1. Initialization

Creating a first/starter solution, no matter how simple/bad it is.

2. Iterate

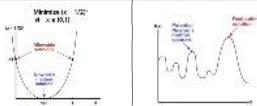
Report \rightarrow Find an improving direction t
 • Using a step size α to move along it
 • New solution = old solution + αt

3. Stop

...when solution doesn't change much, or time runs out.

Algorithm from Calculus
Example { Newton's method: finding root of $f(x)$
 • Current solution x_n at step n
 • $f'(x_n) = -f(x_n)/\Delta x$
 • Step size = Δx
 • Improving direction = $-f'(x_n)$

Convex Optimization Problem	Non-Convex Optimization Problem
Guaranteed to find optimal solution	NOT Guaranteed to find optimal solution
May Converge to Infeasible Solution	
May find only a local optimum	



Non-Parametric Models

Typical hypothesis testing assume that we know the underlying distribution we're testing against.

But that's not always the case. Sometimes all we have is a response and we don't have a good model for the distribution.

Non-parametric tests: use when nothing is known about underlying distribution

Distribution-free data, but we need just enough data to figure out where the middle value is

McNemar's Test

Used for comparing results on pairs of responses/data points where 2 approaches were used on the same thing.

- Does the following:
1. Throw out all the cases where the results are the same
 2. Test (anti-distributed) responses using the binomial distribution to see whether we'd expect results this extreme, or more extreme by luck.

Two competing treatments for a virus:
 • A: successful on 61/100
 • B: successful on 68/100

Scenarios 1:
 32 cases: neither worked
 61 cases: both A and B worked
 7 cases: B worked, A did not
 0 cases: A worked, B did not
 • Conclude that B is not better
 $p=0.01$ (reject)

Scenarios 2:
 12 cases: neither worked
 41 cases: both A and B worked
 27 cases: B worked, A did not
 2 cases: A worked, B did not
 • Conclude that B is not better
 $p=0.08$ (reject)

McNemar's (binomial) test: only consider where A and B are different

SUMMARY

Non-Parametric Test	Assumptions	How It's Done Other Notable Attributes of Test Uses	Other Notable Attributes of Test Uses	Uses
---------------------	-------------	---	---------------------------------------	------

Scenario 1	32 cases: neither worked 67 cases: both A and B worked 7 cases: B worked, A did not 6 cases: A worked, B did not + Conclude that B is better p=0.01 (reject)
Scenario 2	12 cases: neither worked 41 cases: both A and B worked 22 cases: B worked, A did not 20 cases: A worked, B did not + Conclude that B is not better p=0.38 (reject)

McNemar's (chi-squared) test: only consider where A and B are different

To use McNemar's test, we don't need to know anything about the distribution of how much is used.
It's just a comparison of a pair of results.

Wilcoxon Signed Rank Test for Medians

Answers the following questions:

- Is the median of the distribution different from m ?
- Do two sets of paired samples have the same median?

Does the following:

1. Rank $\{y_1 - m\}, \{y_2 - m\}$ from smallest to largest
2. $W = \sum_{y_i > m} rank(y_i - m)$ = sum of all ranks where $y_i > m$
3. Get p-value test for W

→ Answering whether ranks indicate that the two sets of observations are likely to be different

Mann-Whitney Test

Used for two data set analysis, but not paired samples.

Does the following:

1. Rank all observations together
2. $U = \sum_{x_i < y_j} rank(x_i) + \frac{n(n+1)}{2}$
3. Find significance of U (need software or a table)

Summary:

Use parametric tests when nothing is known about its underlying distribution.

Two Data Sets:

- McNemar's Test (paired yes/no)
- Wilcoxon Signed Rank Test (paired numeric data)
- Mann-Whitney (unpaired)

One Data Set:

- Wilcoxon Signed Rank Test (compare possible median)

SUMMARY

Non-Parametric Test	Assumptions	How It's Done Other Notable Attributes of Test Uses	Other Notable Attributes of Test Uses	Uses
McNemar's Test	None needed. This method doesn't care about the underlying distribution	<ul style="list-style-type: none"> 1. Throws out all the cases where the results are the same. 2. Tests the responses (NOT the underlying distribution) using the binomial distribution to see whether we'd expect results this extreme or more extreme just by luck 	We don't need to know anything about the distribution. It's just a comparison of the results	Used for comparing results on pairs of responses, data points where two different approaches were used on the same thing.
Wilcoxon Signed Rank Test for Medians	Distribution is continuous and symmetric	<p>Median difference Test Given responses y_1, \dots, y_n:</p> <ol style="list-style-type: none"> 1. Rank $y_1 - m , \dots, y_n - m$ from smallest to largest 2. $W = \sum_{y_i > m} rank(y_i - m)$ only for y_i values where $y_i > m$ 3. Get p-value test for W <p>Sample Median Equivalence Test (paired test) Given pairs $(y_1, z_1), \dots, (y_m, z_m)$ from observations y & z, use $y_1 - z_1 , \dots, y_m - z_m$ for rank test</p>	Calculating p-values here are harder than McNemar's Test. The Wilcoxon test is like a normal distribution test	Answers the questions: 1. Is the median of the distribution different from m ? 2. Do two sets of paired samples have the same median?
Mann-Whitney Test	Samples are independent from each other	<p>Given independent observations y_1, \dots, y_n and z_1, \dots, z_m:</p> <ol style="list-style-type: none"> 1. Rank all observations together: $y_1, \dots, y_n, z_1, \dots, z_m$ 2. $U =$ smaller of two adjusted rank sums: $U = \min(U_y, U_z)$ where: $U_y = \sum_{j=1}^m rank(y_j) - \frac{n(n+1)}{2}$ $U_z = \sum_{j=1}^n rank(z_j) - \frac{m(m+1)}{2}$ 3. Find significance of U (need software or table) 	—	Used for Two Data set analysis, but not paired samples

Bayesian Models

Based on conditional probability, aka Bayes' Theorem

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

"The probability of A given the probability of B given A, times the probability of A, divided by the probability of B."

The problem with Bayes' Theorem is that it can feel counter-intuitive.

Example: Medical test for a disease

- True positives: 98%
- False positives: 8%
- 1% of population really has disease
- 8.9% of people test positive
- If someone tests positive, what is the probability they have the disease?

• A: has the disease
 • B: tested positive
 • $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$
 • After testing positive
 • a person has only an 11% chance of having the disease
 Why?
 So many more people don't have the disease...
 ...so there are many more false positives than true positives.

Empirical Bayes Modeling

Useful when an overall distribution of something is known or estimated, but there's only a little data available for a specific case.

Example:

- predicting outcomes of the NCAA basketball tournament
- Difference X in points scored by the home team and road team approximately normal: $X \sim N(m + b, \sigma^2)$
- b : home court advantage
- m : true difference in the teams' strength (unknown)
- σ^2 : variance
- The differences between college basketball teams' strength $\sim m - N(0, \tau^2)$

→ Observed data
 • x = observed point difference in game
 • m = real difference between the two teams, $m \neq x$
 • $P(M = m)|X = x = \frac{P(X=x|m+b, \sigma^2)P(m|m)}{P(X=x)}$
 • $M|X = x \sim N\left(\frac{x^2}{\tau^2 + 1} + b, \frac{\tau^2 \sigma^2}{\tau^2 + 1}\right)$
 • Get σ, τ, b from the past data
 • $P(\text{home team better} | X = x) = \int_0^\infty P(M = m|x) dm$

→ Example: home team won by 20 points
 • Estimate home court advantage $b = 4$ points
 • Standard deviation in team strength $\tau = 6$ points
 • Off the top of my head:
 • About 4 points due to home court
 • About 12.5 points due to random variation
 • Only 2.5 points due to the difference between teams
 • A lot more variance due to randomness

In Summary:

- We can:
- 1) Take a single observation,
- 2) Combine with a broader set of observations that are relevant in general (but not specifically);
- 3) Then make a deduction/prediction.

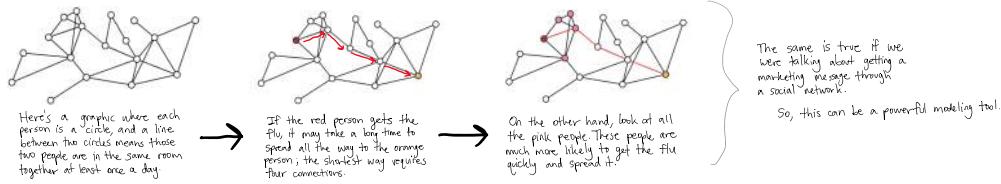
So, Bayesian models work especially in the absence of lots of data.

Communities

An analysis of large, interconnected populations may have impacts on:

- Marketing messages spread through social media networks
- Ideas, cultures, trends through migrations
- Computer viruses spread through a network
- Ideas, news, trends spread from template to language over time
- Information spreading communities
- and more

A look into highly-connected subpopulations...



Some Definitions....

Community: a set of circles that's highly connected within itself

Graph:
 Circles = nodes/vertices
 Lines = arcs/edges
 Clique = a set of nodes that all have edges between each other

Louvain Algorithm: heuristic to decomposing a graph into communities

Maximize the modularity of a graph

- a_{ij} : weight on the arc between nodes i and j
- w_i : total weight of arcs connected to i
- W : total weight of all the arcs
- Modularity = $\frac{1}{2W} \sum_{i,j} \text{in same community } (a_{ij} - \frac{w_i w_j}{2W})$

Steps:

Step 0

Each node is its own community



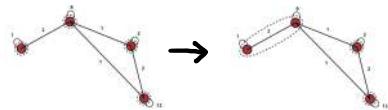
Step 1

Repeat...
 Make biggest modularity increase
 by moving a node from its
 community to an adjacent
 node's community
 ...until no move increases modularity



Step 2

Each community is a super-node
 Repeat Step 1 using super-nodes



Basically, whichever move creates the biggest modularity increase, that's where node i goes. Or if nothing increases modularity, node i just stays where it is. It does that for every node and then starts again with the first node and goes through them all again until there's no way to improve modularity by moving a node.

Neural Networks & Deep Learning

Training system to react to patterns human brains can react to, even if we don't understand the pattern ourselves.

Used in:

- Natural Language Processing
- Speech recognition
- Image Recognition

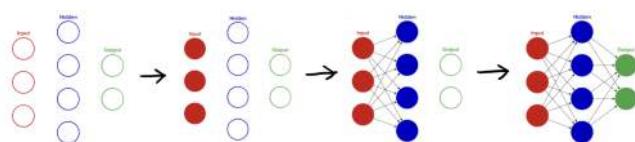
3 levels of neurons:

- Input Level
- Hidden Level
- Output Level

Each neuron:

- Gets inputs from previous layer
- Calculates function of weighted inputs
- Gives its output to next layer

Deep Learning = many deep layers



Competitive Models

Descriptive models: get an understanding of reality

Predictive models: find hidden relationships; predict the future

Prescriptive models: find the best thing to do

Assumes the system
 does not react

...But what if the system reacts, and it reacts intelligently?



→ Use analytics to consider all sides to the system.

Game Theory: competitive decision-making

Cooperative Game Theory: competition + cooperation

Need to think of varying strategies:

- counter-strategies for competitor
- counter-counter strategy (iterative)
- Sequential game: see other side's strategy before making your decision (via user)
- Pure strategy: just do one thing
- Mixed strategy: randomized strategy (according to probabilities)

↳ i.e. rock paper-scissors

Information {

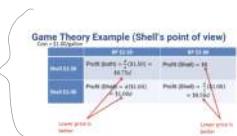
- Perfect information: know all about everyone else's situation; balanced/symmetric, i.e. chess
- Imperfect information: some have more information than others; unbalanced/asymmetric, i.e. competing restaurants

Sum {

- Zero-sum: whatever one side gets, the other side loses, i.e. 1 or -1 loss like rock paper-scissors
- Non-zero-sum games: total benefit might be higher or lower, i.e. economy

Game Theory Example

- Two gas stations  
- Set price: \$2.50 or \$2.00
- Same price: 50/50 demand
- Otherwise: all demand to lower priced one
- What's the best price to choose?



Stable Equilibrium

- Neither station has incentive to change
- In this case...
 - Outcome is worse than if they had agreed to both charge the higher price
 - "Prisoner's dilemma"
 - Even if agree on higher price, both sides have incentive to cheat
 - So both end up with lower price
- Not always true in game theory



Same Example with Smaller Profit Margin



Extension of Example

- Choose any price they want
- BP price = P_B
 - Shell: If $P_B > P_S$, then profit = \$0 (will switch)
 - If $P_B = P_S$, then profit = $(P_S - C) \cdot D$
 - If $P_B < P_S$, then profit = $(P_B - C) \cdot D$
- So, Shell might price slightly lower than BP
- Then, BP might price slightly lower than Shell
 - Then, Shell might price slightly lower than BP
- Both keep lowering prices until price is about equal to the cost

Key Notes

Saturday, April 30, 2022 2:35 PM

Maybe small sample size...
↑ factors, ↓ data → use ridge regression

Maybe linear regression fits data too well.
Ridge regression finds a new line that doesn't fit the training data as well.

(Least Squares)
Linear Regression: no bias \Rightarrow Lut, ↑ variance

Ridge Regression: ↑ bias \Rightarrow so, ↓ variance

Lasso Regression can exclude useless variables \rightarrow good when so many factors

Ridge Regression more useful if we know all the factors are important.

Elastic-Net Regression good if we don't know the situation since $EL = \text{Ridge} + \text{Lasso}$ (best of both worlds)

Classification tree \rightarrow true/false (discrete) leaves

Regression tree \rightarrow numeric data leaves

Logistic Regression \rightarrow discrete values (yes/no) for axis (y-axis namely)

SVM: maximize gap 
(margin)